

¹
²
**Interactive Linear Algebra with
R and Python**

³
SAMUEL S.P. SHEN

⁴ Department of Mathematics and Statistics, San Diego State University,
⁵ and
⁶ Scripps Institution of Oceanography, University of California, San Diego

⁷ August 2024 Version for SDSU Course Math 524

Contents

| | | |
|--|-------------|------|
| <i>Preface</i> | <i>page</i> | viii |
| <i>Acknowledgements</i> | xiii | |
| 1 Basic Matrix Operations | 1 | |
| 1.1 Matrix as a data array | 1 | |
| 1.2 Matrix algebra | 3 | |
| 1.2.1 Matrix equality, addition and subtraction | 3 | |
| 1.2.2 Matrix multiplication | 4 | |
| 1.3 Some Useful Operations for Matrices in Data Science but Not Covered in Traditional Linear Algebra Books | 10 | |
| 1.3.1 Delete a Row or Column or Both; a Sub-Matrix | 10 | |
| 1.3.2 Insert Rows or/and Columns to a Matrix | 11 | |
| 1.3.3 Statistics of the row or column data | 12 | |
| 1.3.4 Sweep a matrix by a vector | 12 | |
| 1.3.5 Conversions between a Vector and a Matrix | 14 | |
| 1.3.6 Reduce Dimensions of an n-Dimensional Array | 14 | |
| 1.4 A Set of Linear Equations | 15 | |
| 1.5 Eigenvalues and eigenvectors of a square space matrix | 16 | |
| 1.5.1 Matrices of data anomalies, standardized anomalies, covariance, and correlation | 17 | |
| 1.5.2 Eigenvectors and their corresponding eigenvalues | 19 | |
| 1.6 An SVD representation model for space-time data | 20 | |
| 1.6.1 Space-Time Data Matrix and Its decomposition | 20 | |
| 1.6.2 SVD of the Spae-Time Anomaly Data Matrix and Eigen- value Problem of a Covariance Matrix | 22 | |
| 1.7 Image analysis using SVD and R | 25 | |
| 1.7.1 Data for a color photo and its grayscale figure | 25 | |
| 1.7.2 SVD of the grayscale photo data | 26 | |
| 1.7.3 Reconstructing a color photo from the data of monotone photos | 29 | |
| 1.8 Mass balance for chemical equations in marine chemistry | 31 | |
| 1.9 Multivariate linear regression using matrix notations | 32 | |
| 1.10 Chapter summary | 34 | |
| <i>References and Further Readings</i> | 37 | |
| Exercises | 38 | |

| | | |
|----|---|-----|
| 44 | 2 Matrix Theory and Visualization | 41 |
| 45 | 2.1 Matrix Definitions | 41 |
| 46 | 2.2 Fundamental Properties of Matrices | 44 |
| 47 | 2.3 Some basic concepts and theories of linear algebra | 49 |
| 48 | 2.3.1 Linear equations | 49 |
| 49 | 2.3.2 Linear transformations | 50 |
| 50 | 2.3.3 Linear independence | 51 |
| 51 | 2.3.4 Determinants | 52 |
| 52 | 2.3.5 Rank of a matrix | 53 |
| 53 | 2.4 Eigenvectors and eigenvalues | 56 |
| 54 | 2.4.1 Definition of eigenvectors and eigenvalues | 56 |
| 55 | 2.4.2 Properties of eigenvectors and eigenvalues for a symmetric matrix | 60 |
| 56 | 2.5 Hadamard and Other Matrix Multiplications | 62 |
| 58 | 2.5.1 Hadamard Product of Two Matrices of the Same Dimensions | 63 |
| 59 | 2.5.2 Jordan Product of Two Matrices of the Same Dimensions | 63 |
| 60 | 2.5.3 Commutator of Two Matrices of the Same Dimensions | 63 |
| 61 | 2.5.4 Outer Product of Two Vectors and Two Matrices | 64 |
| 62 | 2.5.5 Kronecker Product of Two Matrices of the Same Dimensions | 66 |
| 63 | 2.6 Direct Sum of Two Matrices | 67 |
| 64 | 2.7 Visualization of Eigenvalues and Eigenvectors for a Covariance Matrix | 67 |
| 66 | 2.8 Singular Value Decomposition | 68 |
| 67 | 2.8.1 SVD formula and a simple SVD example | 68 |
| 68 | 2.9 SVD for the standardized sea level pressure data of Tahiti and Darwin | 74 |
| 69 | 2.10 Chapter summary | 76 |
| 71 | <i>References and Further Readings</i> | 78 |
| 72 | Exercises | 79 |
| 73 | 3 Matrix Applications to Machine Learning | 83 |
| 74 | 3.1 K-means clustering | 83 |
| 75 | 3.1.1 K-means setup and trivial examples | 84 |
| 76 | 3.1.2 The number of exhaustive K clusters from N points | 90 |
| 77 | 3.1.3 A K-means algorithm | 91 |
| 78 | 3.1.4 K-means clustering for the daily Miami weather data | 93 |
| 79 | 3.2 Support vector machine | 103 |
| 80 | 3.2.1 SVM for a system of three points labeled in two categories | 108 |
| 81 | 3.2.2 SVM mathematical formulation for a system of many points in two categories | 112 |
| 82 | 3.3 Random forest method for classification and regression | 118 |
| 83 | 3.3.1 RF flower classification for a benchmark iris dataset | 118 |
| 84 | 3.3.2 RF regression for the daily ozone data of New York City | 125 |
| 85 | 3.3.3 What does a decision tree look like? | 129 |

| | | |
|-----|--|-----|
| 87 | 3.4 Neural network and deep learning | 132 |
| 88 | 3.4.1 An NN model for an automated decision system | 132 |
| 89 | 3.4.2 An NN prediction of iris species | 139 |
| 90 | 3.5 Chapter summary | 143 |
| 91 | <i>References and Further Readings</i> | 144 |
| 92 | Exercises | 145 |
| 93 | 4 Matrix Applications to Regression Models | 148 |
| 94 | 4.1 Simple linear regression | 148 |
| 95 | 4.1.1 Temperature lapse rate and an approximately linear model | 148 |
| 96 | 4.1.2 Assumptions and formula derivations of the single variate | |
| 97 | linear regression | 152 |
| 98 | 4.1.3 Statistics of slope and intercept: Distributions, confidence | |
| 99 | intervals, and inference | 163 |
| 100 | 4.2 Multiple linear regression | 176 |
| 101 | 4.2.1 Calculating the Colorado TLR when taking location | |
| 102 | coordinates into account | 176 |
| 103 | 4.2.2 Formulas for estimating parameters in the multiple linear | |
| 104 | regression | 179 |
| 105 | 4.3 Nonlinear fittings using the multiple linear regression | 181 |
| 106 | 4.3.1 Diagnostics of linear regression: An example of global | |
| 107 | temperature | 181 |
| 108 | 4.3.2 Fit a third order polynomial | 186 |
| 109 | 4.4 Linear Regression by Weighted Least Squares | 189 |
| 110 | 4.5 Chapter summary | 190 |
| 111 | <i>References and Further Readings</i> | 191 |
| 112 | Exercises | 192 |
| 113 | Appendix A A Tutorial of R and RStudio | 195 |
| 114 | A.1 Download and install R and R-Studio | 195 |
| 115 | A.2 R Tutorial | 196 |
| 116 | A.2.1 R as a smart calculator | 197 |
| 117 | A.2.2 Define a sequence in R | 198 |
| 118 | A.2.3 Define a function in R | 199 |
| 119 | A.2.4 Plot with R | 199 |
| 120 | A.2.5 Symbolic calculations by R | 200 |
| 121 | A.2.6 Vectors and matrices | 201 |
| 122 | A.2.7 Simple statistics by R | 204 |
| 123 | A.3 Online Tutorials | 205 |
| 124 | A.3.1 YouTube tutorial: for true beginners | 206 |
| 125 | A.3.2 YouTube tutorial: for some basic statistical summaries | 206 |
| 126 | A.3.3 YouTube tutorial: Input data by reading a csv file into R | 206 |
| 127 | A.4 Chapter summary | 208 |

| | | |
|-----|---|-----|
| 128 | <i>References and Further Readings</i> | 210 |
| 129 | Exercises | 211 |
| 130 | Appendix B Visualization of Matrices | 214 |
| 131 | <i>References and Further Readings</i> | 246 |
| 132 | <i>Index</i> | 251 |

¹³³ Give the pupils something to do, not something to learn; and doing is
¹³⁴ of such a nature as to demand thinking; learning naturally results.

¹³⁵

¹³⁶ — John Dewey

Preface

138 According to [Encyclopedia.com](#), “linear algebra originated as the study of linear
139 equations.” Linear algebra deals with vectors, matrices and vector spaces. Before
140 the 1950s, it was part of Abstract Algebra (Tucker 1993). For example, the book
141 “A Survey of Modern Algebra” by Garrett Birkhoff and Saunders Mac Lane (3rd
142 edition, Macmillan Company, New York, 1965) contains three chapters of linear
143 algebra materials out of its 15 chapters: Chapter 7 (Vectors and Vector Spaces),
144 Chapter 8 (The Algebra of Matrices), and Chapter 10 (Determinants and Canonical
145 Forms).

146 In 1965 the Committee on the Undergraduate Program in Mathematics, Mathe-
147 matical Association of America (MAA), outlined the following topics for a stand-
148 alone linear algebra course: Linear systems, matrices, vectors, linear transfor-
149 mations, unitary geometry with characteristic values.

150 The first stand-alone book that was named “Linear Algebra” in the United States
151 might be that by Charles W. Curtis of the University of Wisconsin-Madison, pub-
152 lished in 1963 by Allyn and Bacon, Inc., Boston. The book title is “Linear Algebra:
153 An Introductory Approach.” This book contains all the materials defined by MAA
154 in 1965. Most linear algebra books today would also cover these materials, plus some
155 modern materials on singular value decomposition (SVD) (Strang 2016). Modern
156 advanced materials may include spectral theorem, Jordan forms, polynomials, and
157 Cayley-Hamilton theorem (Garcia and Horn 2022).

158 Paul R. Halmos’ book “Finite-Dimensional Vector Space” was mostly about lin-
159 ear algebra although did not name so. This book was first published in 1942 and had
160 its second edition in 1958. However, the book missed an important part of linear
161 algebra: The solution of a set of linear equations, which became extremely impor-
162 tant in the 1960s when computer applications for solving many linear equations in
163 engineering were popular due to the invention of the finite element method.

164 Although Jean Dieudonne’s book “Linear Algebra and Geometry” in 1964 bore
165 the name of Linear Algebra, focused mainly on linear transforms, and had no solu-
166 tions of linear equations either. Ross A. Beaumont’s 1965 book also had its name
167 “Linear Algebra”, but again contained no materials on linear equations.

168 Modern linear algebra books have more materials on rectangular matrices, SVD,
169 data science applications, graphic visualization, and computer codes. Our book is
170 in this category in terms of materials. In the methodology aspect of learning and
171 teaching, we use progressive education pedagogy and emphasize intuition, graphics,
172 storytelling, computing, and rigorous proof.

173 What are the distinguished features of this book?

174 Every chapter starts with elementary materials and common sense, and can be
175 understood without reading previous chapters. Almost all the chapters include ad-
176 vanced materials which are prepared for the second-semester linear algebra course,
177 or called advanced linear algebra. These advanced materials can be used as a re-
178 search handbook or as advanced examples to support a research project. Different
179 from the traditional linear algebra book, chapters are entitled according to the na-
180 ture of the relevant theory, e.g., linear transform and vector space, our chapters
181 are entitled according to the usage or functions of a method, e.g., Data matrices,
182 operations, and visualizations for Chapter 1, and Space-time decomposition: SVD
183 for Chapter 2. We follow the RUM principle when selecting materials. Here R is for
184 Relevant to both courses and research in the science; U is for Useful to not only the
185 courses and research, but also to future jobs; and M is for Modern, incorporating
186 recent progress and tools, especially progress in computing and data science. In this
187 book, our modern materials are reflected in SVD applications, matrix visualization,
188 machine learning, and the use ChatGPT, Python, and R.

189 Meeting these RUM criteria in courses, including linear algebra, will require time
190 and commitment to changing the curricula of relevant study programs to better
191 support students learning. Progress will be incremental. Mathematics courses and
192 instructors can both evolve to help meet the objective of well-educated students
193 prepared for successful careers.

194 In terms of the book style, we follow the principle of interaction. We make sure
195 that each method has a computer code to allow a reader to interact with the book.
196 Traditional linear algebra books, like most other math books, often stack with
197 definitions of mathematical objects, such as subspace, and isomorphism, dense for-
198 mulas, and awkwardly-worded theorems. The theorem statements are only compre-
199 hensible by professional mathematicians in the name of “mathematical rigor,” with
200 which mathematicians are extremely proud of.

201 Bearing in mind that most students of linear algebra are not going to become
202 professional mathematicians, this book tries to do it differently. We start with com-
203 mon sense, data, visualization, examples, and interpretations that lead to names,
204 definitions, and theorems. Most math objects in our book are named after the ob-
205 jects have appeared in a computational process or an algorithm, so that students
206 have an image in their mind before we name them. This is like that the 10-year-old
207 Mike learns to name Uncle Tom after seeing him, in contrast to the legendary Uncle
208 Tom who suddenly appears in front of Mike.

209 Despite our emphasis on the details of modern and useful materials, the book still
210 includes almost all the contents of a conventional linear algebra book, such as vector
211 spaces, inner product space, linear equations, orthogonality, eigenvectors, SVD,
212 linear transform, polynomials, operators, spectral theory, Jordan forms, Cayley-
213 Hamilton theorem, and determinants. Of course, some materials are very brief,
214 such as determinants.

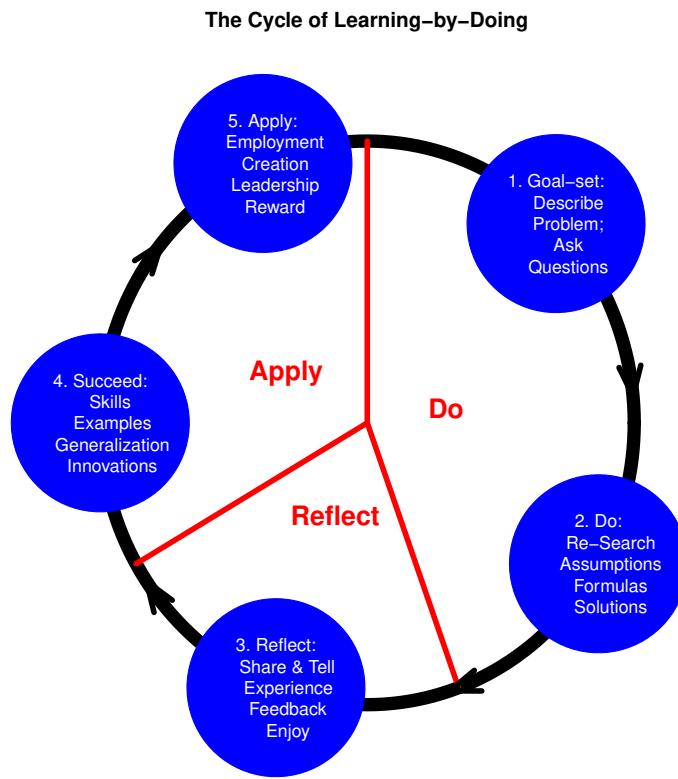
215 Another feature of this book is to minimize the use of mathematical symbols
216 and notations. We understand that professional mathematicians do not mind fancy

217 notations, but non-mathematicians may be scared away by them. Instead, we use
218 some repeated notations for matrices. When we introduce a name and a notation,
219 we attempt to justify its use from the perspectives of history, interpretation, and
220 applications. This way may help a reader to develop a picture in her mind when
221 reading a method, or a definition. We use visualization as much as we can, such
222 as the visualization of a data matrix, and that of left and right eigenvectors from
223 SVD.

224 **What is the philosophy of our pedagogy?**

225 We follow the education theory of *learning-by-doing*, which in our case means
226 using your computer and our R or Python code to interact with our book and
227 modifying our computer codes and websites to analyze your own data and generate
228 corresponding figures. Learning-by-doing is the core methodology of the progres-
229 sive education of John Dewey (1859-1952), an American philosopher and educator.
230 Dewey felt that the experience of students and teachers together yields extra value
231 for both. Instructors are not just to lecture and project authority, instead they
232 are to collaborate with students and guide to students to gain experience of solv-
233 ing problems of their interest. Although Dewey's education theory was established
234 initially for school children, we feel that the same is applicable to undergraduate
235 and graduate students. Our way of learning-by-doing is to enable students to use
236 R or Python code and other resources in the book and its website to reproduce
237 the figures and numerical results in the book. Further, students can modify the
238 computer code and solve their own problems, such as visualizing the health data or
239 weather data in a similar format and of their own interest, or analyzing their own
240 data using a similar or a modified method. Thus, the audience interaction with the
241 book is the main innovative feature of our book, which allows the audience to gain
242 experience of practicing, thinking, applying, and consequently understanding. The
243 ancient Chinese educator Confucius (551-479 BC) said that "*I hear, and I forget; I see, and I remember; and I do, and I understand.*" Although John Dewey and
244 Confucius were approximately 2,500 years apart, they shared a similar philosophy
245 of learning-by-doing.

247 As illustrated in Fig. 0.1, our pedagogy has three stages: Do, reflect, and apply.
248 The author has systematically practiced this pedagogy since 2015 when he created
249 and taught the course *Climate Mathematics* at the Scripps Institution of Oceanog-
250 raphy. Usually, he presents a question or a problem at the beginning of a class.
251 Then he asks students to orally ask the same question, or describe the same prob-
252 lem, or answer his question in their own words. For example, why the tiny amount
253 of carbon dioxide in the atmosphere is important for climate change? Next, he and
254 students search and re-search for data and literature, work on the observed carbon
255 dioxide data at Mauna Loa using computer data visualization to see the 30% in-
256 crease of carbon dioxide in the atmosphere since 1958, and discuss the structure
257 of greenhouse gasses whose molecules have three or more atoms. Next, he encour-
258 ages his students to share their experiences with their grandparents, other family

**Fig. 0.1**

John Dewey's pedagogy of learning-by-doing: A cycle from a problem to skill training via solving the problem, to reflection, to success and applications, and to reward and new problems. It is a cycle of understanding and progress. The size of each piece of the "Do, Reflect, Apply" pie approximates the proportion of the learning effort. "Do" is emphasized here, indicated by the largest piece of the pie.

members, or friends. Finally, students apply the skills gained to solve their own problems with their own data, by doing homework, working on projects, finding employment, or making innovations. In this cycle, students have gathered experience and skills to improve their life and to engage with the community. In the short term, students trained in this progressive cycle of learning-by-doing have a better chance to become a good problem solver, a smooth story narrator, and an active leader in a research project in a lab or an internship company, and consequently become competitive in the job market. In the long term, the students trained in this cycle are likely to become a life-time learner and educator. John Dewey said that "Education is not preparation for life but life itself." We would like to modify this to that "Education is not only preparation for a better life, but also is life itself." Dewey's progressive education theory is in a sharp contrast to the traditional

learning process based on the logic method, which aims at cultivating highly achieved scholars. The commonly used pedagogy of lecture-read-homework-exam often uses the logic-based approach. The climax of the logic-based education is that the instructor has the pleasure to enjoy presenting her method and theory, while students are so creative that they will produce a new or a better theory. Many outstanding scholars went through their education this way, and many excellent textbooks were written in this approach. However, our book does not follow this approach, since our book is written for the general population of students, not just for the elite class or even scholars-to-be. If you wish to be such an ambitious scholar, then you may use our book very differently: you can read the book quickly and critically for gaining knowledge instead of skills, and skip our reader-book interaction functions.

Our pedagogy is result-oriented. Our philosophy is to minimize the mathematical challenge, but to deepen the understanding of the ideas using visual tools and using story-telling experience. Our audience will be able to make an accurate problem statement, ask pointed questions, set up linear algebra models, solve the models or establish theorems, prove the theorems in both intuitive and rigorous ways, and interpret solutions and theorems using both daily life language and rigorous mathematical language. Therefore, instead of training the few “mechanics” with a high risk of failure, we wish to train a large number of good “drivers” with a large probability of success, although this book also tries to inspire a few students to become “mechanics.”

292 Samuel S. P. Shen, San Diego, California, August 2024

293 References:

- 294 1. Beaumont, R. A., 1965. Linear Algebra. Harcourt, Brace & World, Inc., New
295 York.

296 2. Birkhoff, G. and Mac Lane, S., 2017. A survey of modern algebra. 3rd edition,
297 Macmillan Company, New York, 1965.

298 3. Curtis, C.W., 1963. Linear algebra: an introductory approach. Allyn and Bacon,
299 Inc., Boston.

300 4. Garcia, S.R. and Horn, R.A., 2023. Matrix Mathematics: A Second Course in
301 Linear Algebra. 2nd edition, Cambridge University Press.

302 5. Grothendieck, A., Dieudonn, J. and Dieudonn, J., 1964. Elments de gomtrie
303 algrique. Hermann, Paris.

304 6. Halmos, P.R., 1958. Finite-dimensional vector spaces. 2nd edition, D. Van Nos-
305 trand Company Inc., Princeton.

306 7. Strang, G., 2016. Introduction to linear algebra. 5th edition, Wellesley-Cambridge
307 Press.

308 8. Tucker, A., 1993. The growing importance of linear algebra in undergraduate
309 mathematics. The college mathematics journal, 24(1), pp.3-9.

Acknowledgements

³¹¹ We thank the friendly editorial team of Cambridge University Press (CUP) for their
³¹² professional and high quality work.

316

317 Data matrices appear everywhere in our daily life, such as the daily body weight
318 data of 100 clients of your diet clinic in the last 21 days. If you put each client's
319 data on a row, then the row has 21 data entries. This row may be called a data
320 sequence, or a row vector. You stack the data sequence of all the 100 clients together
321 one above another ordered according to their last name. You thus have created a
322 rectangular 2-dimensional array of numbers, a 100-by-21 matrix. This data matrix
323 has 100 rows and 21 columns.

324 You may add client names for the row names and date for the column names
325 to the data matrix and form a data frame. Data frame is a commonly used concept
326 in R and Python data analysis. Data frame has a nice presentation of data and
327 also allows you to name the data in rows and columns, such as John Smith's body
328 weight data, in contrast to referring to the 83rd row.

329 This chapter is limited to (i) Describing the basic matrix operations commonly
330 used in your coursework or career, e.g., eigenvector decomposition of a matrix and
331 solution of linear equations; (ii) presenting matrix application examples of data
332 science, e.g., sub-matrix generation, array-matrix conversion, and data statistics;
333 and (iii) solutions of linear equations.

334

1.1 Matrix as a data array

335

336 In general, a matrix is a rectangular array of numbers or symbols or expressions,
337 which are called elements, arranged in rows or columns. A table such as that shown
338 in Fig. 1.1 is a matrix, consisting of N rows and Y columns of numbers. Figure
339 1.1 shows the 10 ten-year annual precipitation anomalies from the year 1900 to
340 1909 for the 15 five-degree latitude-longitude boxes centered at $2.5^{\circ}E$ longitude for
341 different latitudes over the Northern Hemisphere ranging from latitudes $2.5^{\circ}N$ to
342 $72.5^{\circ}N$. The matrix shown in Fig. 1.1 thus has 15 rows ($N = 15$), and 10 columns
343 ($Y = 10$). An anomaly of a climate parameter is defined as its actual value minus
344 its normal value that is an average of 30 or more years.

345 Many other types of climate data can also be represented as matrices (which
346 is the plural of matrix). Precipitation data [units: mm/day] at multiple stations
347 and multiple days can also form a matrix, normally with stations [each marked
348 by a station identifier, or station ID] represented in rows, and time [units: day]
349 represented in columns. The daily minimum surface air temperature (Tmin) data

| Lat | Lon | 1900 | 1901 | 1902 | 1903 | 1904 | 1905 | 1906 | 1907 | 1908 | 1909 |
|------|-----|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 2.5 | 2.5 | 0.283240 | -0.131860 | -0.190500 | 0.160040 | -0.878110 | 0.080356 | 0.059193 | -0.136900 | 0.200420 | 0.822600 |
| 7.5 | 2.5 | 0.172670 | 0.830550 | -0.180350 | -0.203630 | -0.238590 | 0.425310 | 0.002805 | 0.102780 | 0.254050 | 0.516200 |
| 12.5 | 2.5 | 0.024392 | 0.152030 | -0.034115 | -0.062696 | -0.192070 | 0.074360 | 0.201970 | -0.011311 | 0.035259 | 0.272010 |
| 17.5 | 2.5 | 0.006780 | 0.066783 | -0.084581 | -0.008636 | -0.038109 | -0.001092 | 0.088250 | 0.011047 | 0.029358 | 0.082329 |
| 22.5 | 2.5 | 0.021162 | 0.079977 | 0.020016 | -0.022142 | -0.027032 | 0.065704 | 0.012937 | -0.003823 | 0.032545 | 0.028636 |
| 27.5 | 2.5 | 0.049846 | 0.057413 | 0.026621 | 0.019914 | -0.002651 | 0.071242 | 0.012837 | 0.001567 | 0.051857 | 0.099650 |
| 32.5 | 2.5 | 0.107740 | 0.143510 | 0.061613 | 0.076137 | 0.147760 | 0.137890 | -0.074612 | 0.110300 | 0.087752 | 0.126920 |
| 37.5 | 2.5 | 0.128250 | 0.211940 | 0.113010 | 0.027472 | 0.183710 | 0.125550 | -0.267500 | 0.215980 | 0.007609 | 0.055573 |
| 42.5 | 2.5 | 0.158490 | 0.800950 | 0.292690 | 0.172930 | 0.272010 | 0.126370 | -0.017183 | 0.184880 | 0.118980 | 0.200520 |
| 47.5 | 2.5 | -0.112800 | 0.243130 | -0.121630 | -0.076247 | -0.047231 | 0.110160 | 0.080978 | -0.091371 | 0.016172 | -0.060487 |
| 52.5 | 2.5 | -0.199840 | -0.381070 | -0.217570 | -0.107760 | -0.124700 | -0.117470 | -0.062448 | -0.171070 | -0.277650 | -0.132690 |
| 57.5 | 2.5 | -0.076619 | -0.515070 | 0.005342 | 0.016647 | 0.137820 | 0.038041 | 0.131370 | -0.196490 | -0.132480 | 0.014887 |
| 62.5 | 2.5 | -0.261760 | -0.402600 | 0.137200 | -0.214960 | 0.249210 | 0.147550 | 0.866120 | -0.453910 | -0.026134 | 0.053409 |
| 67.5 | 2.5 | 0.034079 | 0.223610 | 0.314090 | -0.044832 | 0.130470 | 0.201260 | 0.554170 | -0.054434 | 0.185870 | 0.308950 |
| 72.5 | 2.5 | -0.119680 | 0.022949 | 0.004324 | -0.050248 | 0.251330 | -0.233080 | -1.043800 | 0.363850 | -0.315400 | -0.113080 |

Fig. 1.1 Annual precipitation anomalies data of the Northern Hemisphere at longitude $2.5^{\circ}E$ [units: mm/day]. The annual total of the anomalies should be multiplied by 365.

for the same stations and the same days form another matrix. In general, a space-time climate data table always forms a matrix. Conventionally, the spatial locations correspond to the rows, and the time coordinate corresponds to the columns.

Another matrix example, taken from everyday life, is that the ages of members of an audience, sitting in a movie theater in seats arranged in rows and columns, also form a matrix. The weights of these audience members form another matrix. Their bank account balances form still another matrix, and so on. Thus, a matrix is a data table, and extensive mathematical methods have been developed in the 20th century to study matrices. Computer software systems, such as R, have also been developed in recent years that greatly facilitate working with matrices.

This chapter will discuss following topics:

- (i) Matrix algebra of addition, subtraction, multiplication, and division (i.e., inverse matrix);
- (ii) Linear equations;
- (iii) Space-time decomposition, eigenvalues, eigenvectors, and the climate dynamics interpretation of a space-time climate data matrix;
- (iv) A matrix application example in balancing the mass in a chemical reaction equation by solving a set of linear equations; and
- (v) A matrix application in multivariate linear regression.

1.2 Matrix algebra

369

370

371 Matrix algebra is quite different from the algebra for a few scalars of x, y, z as
 372 we learned in high school. For example, matrix multiplication does not have the
 373 commutative property, i.e., matrix A times matrix B is not always the same as
 374 matrix B times matrix A . This section describes a set of rules for doing matrix
 375 algebra.

1.2.1 Matrix equality, addition and subtraction

376

377

378 An $m \times n$ matrix A has m rows and n columns and can be written as

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad (1.1)$$

379 or

$$A_{m \times n} = [a_{ij}], \quad (1.2)$$

380 or simply

$$A = [a_{ij}], \quad (1.3)$$

381 where $a_{ij}, i = 1, \dots, m, j = 1, \dots, n$ are the mn elements of the rectangular matrix
 382 A , and $m \times n$ is often called the size or order of a matrix. We say that A is an m
 383 times n matrix, or an m -by- n matrix. The elements of the matrix shown in Fig. 1.1
 384 are the precipitation anomaly data, $m = 15$, and $n = 10$. Therefore, Fig. 1.1 is a
 385 15 times 10 matrix.

386 Matrix $A = [a_{ij}]$ is equal to matrix $B = [b_{ij}]$ if and only if $a_{ij} = b_{ij}$ for all the
 387 elements. That is, A is identical to B . We can understand this by considering the
 388 two identical photos. A black and white photo is a matrix of pixel brightness values.
 389 Two photos are identical only when the corresponding brightness values of the two
 390 photos are the same. Thus, when considering that a matrix is equal to another,
 391 we may regard the equality as two same-size photos, maps, or climate charts being
 392 identical to one another.

393 Matrix addition is simply the addition of the corresponding elements:

$$A + B = [a_{ij} + b_{ij}]. \quad (1.4)$$

394 Of course, two matrices can be added together only when they have the same size
 395 $m \times n$. If the two matrices represent dimensional data, such as precipitation, then
 396 their units must be the same in order to add corresponding elements. However,
 397 each element in a matrix can represent a different climate parameter. For example,
 398 a climate data matrix for San Diego for 24 hours may have its first row representing
 399 temperature, the second precipitation, the third atmospheric pressure, the fourth

400 relative humidity, etc, while the first columns represent time from 1:00 (i.e., 1:00
 401 AM) to 24:00 (i.e., 12:00 AM).

402 Matrix subtraction is defined in a similar way:

$$A - B = [a_{ij} - b_{ij}]. \quad (1.5)$$

403
 404 **Example 1.1** Matrix subtraction:

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} - \begin{bmatrix} 1 & 2 \\ 3 & 0 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ -2 & -1 \end{bmatrix} \quad (1.6)$$

405
 406 The R code for the above matrix subtraction can be written as follows:

```
408 matrix(c(1,1,1,-1), nrow=2) - matrix(c(1,3,2,0), nrow=2)
409 # [,1] [,2]
410 #[1,] 0 -1
411 #[2,] -2 -1
```

412 One can use matrix subtraction to quantify the differences of climate between
 413 two space-time domains when the climate data for each domain are in a matrix
 414 form.

415 1.2.2 Matrix multiplication ---

416 While the matrix addition and subtraction are similar to those of scalars, the matrix
 417 multiplication has several distinct properties.

419 1.2.2.1 A row vector times a column vector

420 The single-column n-row matrix is often called an n-dimensional vector, or an n-
 421 dimensional column vector. Similarly, one can define an n-dimensional row vector
 422 as a single-row n-column matrix.

423 A row vector of n elements times a column vector of the same number of elements
 424 is equal to a scalar, which is the sum of the products of each pair of corresponding
 425 elements:

$$\begin{bmatrix} a_1 & a_2 & \cdots & a_n \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = a_1b_1 + a_2b_2 + \cdots + a_nb_n \quad (1.7)$$

426 This is also called the dot product of two vectors $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and $\mathbf{b} =$
 427 (b_1, b_2, \dots, b_n) , denoted by

$$\mathbf{a} \cdot \mathbf{b} = a_1b_1 + a_2b_2 + \cdots + a_nb_n. \quad (1.8)$$

428 In 2- or 3-dimensional space, the dot product of two vectors has a simple geo-
 429 metric interpretation:

$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \gamma, \quad (1.9)$$

430 where $\|\mathbf{a}\|$ stands for the length of a vector \mathbf{a} , and γ is the angle between \mathbf{a} and \mathbf{b} .
 431 The proof of the equivalence of the above two expressions (1.7) and (1.9) is given
 432 in Appendix A.

433 When \mathbf{a} is force and \mathbf{b} is the displacement of an object moved by the force, then
 434 $\mathbf{a} \cdot \mathbf{b}$ is the work done by the force on the subject.

435

436 **Example 1.2** Dot product of two vectors in a 2-dimensional space:

$$\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 1 \quad (1.10)$$

437 The same result can be obtained from the geometric formula (1.9). The length of
 438 the first vector is $\sqrt{2}$ since it is the diagonal vector of a unit square in the first
 439 quadrant, and the length of the second vector is 1 since it is the unit square's side
 440 on the x -axis. The angle between the two vectors is 45° . Thus, the dot product of
 441 the two vectors is $\sqrt{2} \times 1 \times \cos(45^\circ) = 1$.

```
442 #R code for dot product
443 #install.packages('geometry')
444 library(geometry)
445 a = c(1, 1)
446 b = c(1, 0)
447 # Calculating dot product using dot()
448 dot(a, b)

449
450 #Another way to compute dot product is to
451 a%*%b
452 #      [,1]
453 #[1,]    1
454
455 #Calculate the angle between two vectors
456 am = norm(a, type="2")
457 am
458 #[1] 1.414214
459 bm = norm(b, type="2")
460 bm
461 #[1] 1
462 angleab = acos(dot(a, b)/(am*bm))*180/pi
463 angleab
464 #[1] 45 degrees
```

465

466

1.2.2.2 A scalar times a matrix

467 A scalar c times a matrix $A = [a_{ij}]$ is formed by multiplying the scalar into every
 468 element of the matrix.

$$c \times A = [c \times a_{ij}]. \quad (1.11)$$

469 The product is again a matrix of the same size. A physically meaningful product
 470 requires that the dimension of the scalar and the dimensions of the matrix elements
 471 are compatible. For example, if the matrix is the price data of many products, and
 472 if the scalar is the number of sales of each product, then the scalar times the matrix
 473 gives the revenue matrix of all the products.

474

475 **Example 1.3** A scalar 3 times a 2-by-2 matrix

$$A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

476 is

$$3 \times A = 3 \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 3 & 3 \\ 3 & -3 \end{bmatrix} \quad (1.12)$$

```
477 #R code for a scalar times a matrix
478 A = matrix(c(1,1,1,-1), ncol = 2, byrow=T)
479 A
480 #      [,1] [,2]
481 #[1,]    1    1
482 #[2,]    1   -1
483 3*A
484 #      [,1] [,2]
485 #[1,]    3    3
486 #[2,]    3   -3
```

487

1.2.2.3 An m-by-n matrix times an n-by-k matrix

488 The multiplication of two matrices is defined as a set of dot products between the
 489 row vectors of the first matrix and the column vectors of the second matrix. Because
 490 of the requirement to form dot products, the column number of the first matrix
 491 must be the same as the row number of the second matrix. $A_{m \times n}B_{n \times k}$ is defined
 492 as a new matrix $C_{m \times k} = [c_{ij}]$, where the element c_{ij} is the dot product of the i th
 493 row vector of $A_{m \times n}$ and j th column vector of $B_{n \times k}$. Therefore, we may write,

$$A_{m \times n}B_{n \times k} = \left[\sum_{l=1}^n a_{il}b_{lj} \right]_{m \times k}. \quad (1.13)$$

495 Thus, $A_{m \times n}B_{n \times k}$ is equal to a matrix of $m \times k$ dot products, and this can be
 496 tedious to compute.

497

498 **Example 1.4** Matrix multiplications:

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} = \begin{bmatrix} 3 & 7 \\ -1 & -1 \end{bmatrix}, \quad (1.14)$$

499 and

$$\begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 4 & -2 \\ 6 & -2 \end{bmatrix}. \quad (1.15)$$

500 In the above formula (1.14), the first element of the right-hand-side matrix
 501 is 3, which is the product of the first row vector of A : (1, 1), and the first column
 502 vector of B : (1, 2). Their dot product is

$$(1, 1) \cdot (1, 2) = 1 \times 1 + 1 \times 2 = 3. \quad (1.16)$$

503 In the same way, one can verify every element of the above multiplication results.

504 These matrix products can be computed by the following computer code.

```
505 #R code for matrix multiplication by command %*%
506 A = matrix(c(1,1,1,-1), nrow=2)
507 B = matrix(c(1,2,3,4), nrow=2)
508 A %*% B
509 #      [,1] [,2]
510 # [1,]    3    7
511 # [2,]   -1   -1
512 B %*% A
513 #      [,1] [,2]
514 # [1,]    4   -2
515 # [2,]    6   -2
```

516 In this example, the second matrix multiplication (1.15) involves the same
 517 matrices as the first one (1.14), but in a different order: If eq. (1.14) is denoted by
 518 AB , then eq. (1.15) is BA . Clearly, the results are different. In general,

$$AB \neq BA \quad (1.17)$$

519 for a matrix multiplication. Thus, matrix multiplication does not have the commu-
 520 tative property which the multiplication of two scalars x and y does have: $xy = yx$.

521

522 1.2.2.4 Transpose matrix, and diagonal, identity and zero 523 matrices

524 Although a space-time climate data matrix often has its rows to mark spatial loca-
 525 tions and columns to mark time, the row and column roles may need to exchange
 526 for some applications, which use rows to mark time and columns to mark spatial
 527 locations. This operation of exchanging rows and columns is called matrix trans-
 528 pose. The transposed matrix of A is denoted by A^t . The columns of A^t are the rows
 529 of A .

530

Example 1.5

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, A^t = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \quad (1.18)$$

531 The transpose can be computed by the following computer code.

```
532 #R code for matrix transpose
533 A = matrix(c(1,2,3,4), ncol = 2, byrow=T)
534 A
535 #[,1] [,2]
536 #[1,] 1 2
537 #[2,] 3 4
538 t(A)
539 #[,1] [,2]
540 #[1,] 1 3
541 #[2,] 2 4
```

542 It is obvious that

$$(A + B)^t = A^t + B^t. \quad (1.19)$$

543 However, a true but less obvious formula is the transpose of a matrix multiplication:

$$(AB)^t = B^t A^t. \quad (1.20)$$

544

545

546 When a matrix whose only non-zero elements are on the diagonal of elements,
547 this matrix is call a diagonal matrix:

$$D = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{bmatrix}, \quad (1.21)$$

548 An identity matrix is a special type of diagonal matrix, whose diagonal elements
549 are all one and whose off-diagonal elements are all zero, and is denoted by I :

$$I = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}, \quad (1.22)$$

550 which plays a role in matrix operations similar to the role of the value 1.0 in the
551 familiar real number system.

552 A zero matrix is a matrix whose elements are all zero. If two matrices are the
553 same, then their difference is a zero matrix.

554 These matrices may be generated by the following computer code.

```

555 #Generate a diagonal matrix
556 D = diag(c(2, 1, -3))
557 round(D, digits = 0)
558 # [,1] [,2] [,3]
559 #[1,] 2 0 0
560 #[2,] 0 1 0
561 #[3,] 0 0 -3
562
563 #Generate an 3-dimensional identity matrix
564 I = diag(3)
565 I
566 # [,1] [,2] [,3]
567 #[1,] 1 0 0
568 #[2,] 0 1 0
569 #[3,] 0 0 1
570
571 #Generate a 2-by-3 zero matrix
572 M = matrix(0, 2, 3)
573 M
574 # [,1] [,2] [,3]
575 #[1,] 0 0 0
576 #[2,] 0 0 0

```

1.2.2.5 Matrix division and inverse

578 The division of a scalar y by another non-zero scalar x can be written as y times
 579 the inverse of x :

$$y/x = y \times x^{-1}. \quad (1.23)$$

580 Thus, the division problem becomes a multiplication problem when the inverse is
 581 found. The inverse of x is defined as $x^{-1} \times x = 1$.

582 Matrix division is defined in the same way:

$$A/B = A \times B^{-1}, \quad (1.24)$$

583 where B^{-1} is the inverse matrix of B defined as

$$B^{-1}B = I, \quad (1.25)$$

584 where I is the identity matrix.

585
 586 **Example 1.6** The R command to find the inverse of matrix A is `solve(A)` as
 587 shown by a numerical example below.

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}^{-1} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & -0.5 \end{bmatrix} \quad (1.26)$$

588 You can compute and verify this “by hand”, or by the following computer code

```

589 #R code to compute the inverse of a matrix
590 A = matrix(c(1,1,1,-1), nrow=2)
591 invA = solve(A)

```

```

592 invA
593 #      [,1] [,2]
594 #[1,]  0.5  0.5
595 #[2,]  0.5 -0.5
596
597 #Verify the inverse
598 invA %*% A
599 #      [,1] [,2]
600 #[1,]  1     0
601 #[2,]  0     1
602

```

603 Finding the inverse matrix of a matrix “by hand” is usually very difficult and
 604 involves a long procedure for a large matrix, say, a 4×4 matrix. Modern climate
 605 models can involve multiple $n \times n$ matrices with n from several hundred to several
 606 million, or even billion. For detailed information, see the excellent text *Introduction*
 607 to *Linear Algebra* by Gilbert Strang (2016).

608 **1.3 Some Useful Operations for Matrices in Data** 609 **Science but Not Covered in Traditional Linear Algebra** 610 **Books**

612 **1.3.1 Delete a Row or Column or Both; a Sub-Matrix**

614 A column of matrix A may contain some unwanted data, or missing data. You wish
 615 to delete this column, say the m th row. The resulting matrix is denoted by $A[-m,]$.
 616 Similarly, you may wish to delete a column. The result is denoted by $A[, -n]$.
 617 Or you may wish to delete the m th row and n th column at the same time. The
 618 result is denoted by $A[-m, -n]$.
 619 Still, another case is that you may wish to choose a sub-matrix, such as from i th
 620 row to the j th row and from k th column to l th column. The result is denoted as
 621 $A[i:j, k:l]$.

622 These operations can be made by the following computer code.

```

623 #R code to delete rows or/and columns
624 A = matrix(1:9, nrow = 3)
625 A
626 #      [,1] [,2] [,3]
627 #[1,]    1    4    7
628 #[2,]    2    5    8
629 #[3,]    3    6    9
630
631 A[-3,]
632 #      [,1] [,2] [,3]
633 #[1,]    1    4    7
634 #[2,]    2    5    8

```

```

635
636 A[, -1]
637 # [,1] [,2]
638 #[1,] 4 7
639 #[2,] 5 8
640 #[3,] 6 9
641 A[-3, -1]
642 # [,1] [,2]
643 #[1,] 4 7
644 #[2,] 5 8
645
646 A[1:2, 2:3] #Sub-matrix
647 # [,1] [,2]
648 #[1,] 4 7
649 #[2,] 5 8
650
651 A[1:2, ] #keep all the columns
652 # [,1] [,2] [,3]
653 #[1,] 1 4 7
654 #[2,] 2 5 8

```

1.3.2 Insert Rows or/and Columns to a Matrix

655 You may wish to insert an extra row or a column to matrix A . The row may be
 656 inserted at any position, on the top, at the bottom, or immediately below the m th
 657 row of A . Below are some examples.

```

658
659
660 #Insert a row or column to a matrix
661 A = matrix(1:4, nrow = 2)
662 br = 5:6
663 bc = 7:8
664 rbind(A, br)
665 # [,1] [,2]
666 #1 3
667 #2 4
668 #br 5 6
669
670 rbind(br, A)
671 # [,1] [,2]
672 #br 5 6
673 # 1 3
674 # 2 4
675
676 rbind(rbind(A[1,], br), A[2,])
677 # [,1] [,2]
678 # 1 3
679 #br 5 6
680 # 2 4
681 Abc = cbind(A, bc)
682 Abc
683 Abc = cbind(A, bc)
684 #
685 # [1,] 1 3 7
686 # [2,] 2 4 8
687

```

```

688 cbind(Abc, A)#stack two matrices
689 #          bc
690 #[1,] 1 3 7 1 3
691 #[2,] 2 4 8 2 4

```

1.3.3 Statistics of the row or column data

For a data matrix, you may wish to calculate the statistical indices of each row or column, such as mean and standard deviation. Below are some examples of these calculations.

```

697 #Row or column statistics
698 A = matrix(1:6, nrow = 2)
699 rowSums(A)
700 #[1] 9 12
701 rowMeans(A)
702 #[1] 3 4
703 rowCumsums(A) #cumulative sum
704 # [,1] [,2] [,3]
705 #[1,] 1 4 9
706 #[2,] 2 6 12
707 colMeans(A)
708 #[1] 1.5 3.5 5.5
709 library(matrixStats) #SD needs the library
710 rowSds(A)
711 #[1] 2 2
712 colSds(A)
713 #[1] 0.7071068 0.7071068 0.7071068

```

1.3.4 Sweep a matrix by a vector

1.3.4.1 Subtraction sweeping

You may wish to subtract a vector of length n from every row of an $m \times n$ matrix. This is an operation of subtraction sweeping.

```

719 #Sweep a matrix by a vector using subtraction: R code
720 A = matrix(1:6, nrow=2, byrow = TRUE)
721 A
722 # [,1] [,2] [,3]
723 #[1,] 1 2 3
724 #[2,] 4 5 6
725 u = 1:3
726 Br = sweep(A, 2, u) #2 means sweep every row
727 Br
728 # [,1] [,2] [,3]
729 #[1,] 0 0 0
730 #[2,] 3 3 3
731 v= 1:2
732 Bc = sweep(A, 1, v) #1 means sweep every column
733 Bc
734 # [,1] [,2] [,3]
735 #[1,] 0 1 2

```

```

736 # [2,]    2    3    4
737
738 c = colMeans(A) #means of each column
739 sweep(A, 2, c) #compute the anomaly data matrix
740 #[1,] -1.5 -1.5 -1.5
741 #[2,]  1.5  1.5  1.5
742
743 sin(A)#function operation on each matrix element
744 #      [,1]     [,2]     [,3]
745 #[1,]  0.8414710  0.9092974  0.1411200
746 #[2,] -0.7568025 -0.9589243 -0.2794155
747
748 A^2 #not equal to A%*%A
749 #      [,1]     [,2]     [,3]
750 #[1,]    1    4    9
751 #[2,]   16   25   36

```

The subtraction sweeping is useful in data science, such as computing the anomalies of data matrix A by subtracting the mean of each column. In this case, you may use R command `c = colMeans(A)`, and the anomaly data matrix would be the result of row sweeping by subtraction.

This example also shows how to apply a function to each element of a matrix, such as R command `sin(A)`. This yields a matrix of elements $\sin(a_{ij})$.

If you wish to sweep by addition, you can just sweep $-u$ or $-v$.

1.3.4.2 Multiplication sweeping

You can also sweep a matrix A by multiplication with a vector v for every column. The R command is `v*A`. This command makes vector $v_{m \times 1}$ time the first n elements of the matrix A , counting from the first column. If $A_{m \times n}$, then $v * A$ makes each column of $A_{m \times n}$ time vector $v_{m \times 1}$ for every pair of corresponding elements. The computer code for an example is as follows.

```

765 #R sweep a matrix by a vector using multiplication
766 A = matrix(1:6, nrow=2, byrow = TRUE)
767 w = 1:2
768 w*A
769 #      [,1]     [,2]     [,3]
770 #[1,]    1    2    3
771 #[2,]    8   10   12
772 A*w # yields the same result as w*A
773 w3 = 1:3
774 #sweep each row by transposing A
775 t(w3*t(A))
776 #      [,1]     [,2]     [,3]
777 #[1,]    1    4    9
778 #[2,]    4   10   18
779 w3*A ##multiplication sweep by row-dimensions not matching
780 #      [,1]     [,2]     [,3]
781 #[1,]    1    6    6
782 #[2,]    8    5   18
783
784 A/w #sweeping by division

```

```

785 # [,1] [,2] [,3]
786 #[1,] 1 2.0 3
787 #[2,] 2 2.5 3

```

When the multiplication sweeping has unmatched dimensions, the R code continues to work, but the sweeping vector goes to the next column to find elements to multiply until its exhaustion, as shown in this computer example `w3*A`.

An application example is to weight the space-time data matrix A by the weight vector w . The weights depend on spatial locations. Then $w*A$ becomes the weighted space-time data matrix, which is often used in climate science, economics, and sociology.

This computer example also shows the sweeping by division, again counting the elements of the matrix from the first column. See R command `A/w` and its output.

1.3.5 Conversions between a Vector and a Matrix

Sometimes data are given in a sequence form, a vector. However, each datum may correspond to specific meanings, such as the month-end body weight data of a group of five people for a year. Each datum corresponds to a specific person for a specific month. You thus may prefer to present the data in a matrix form, with five rows for people and 12 columns for time.

Below are some examples of the conversions between a vector and a matrix.

```

805 #Conversions between a Vector and a Matrix
806 v = c(60, 58, 67, 70, 55, 53)
807 M = matrix(v, nrow = 2) #from vector to matrix
808 M
809 # [,1] [,2] [,3]
810 #[1,] 60 67 55
811 #[2,] 58 70 53
812 c(M) #from matrix to vector by column
813 #[1] 60 58 67 70 55 53
814 c(t(M)) #from matrix to vector by row
815 #[1] 60 67 55 58 70 53

```

1.3.6 Reduce Dimensions of an n-Dimensional Array

Atmospheric temperature data may be naturally regarded as a 4-dimensional array: three dimensions correspond to latitude, longitude, and elevation, and the fourth to time. For some specific space-time data analysis, such as singular value decomposition to be discussed later, we have to convert the 4D array into a 2D matrix, with the rows corresponding to locations and columns to time, i.e., a space-time data matrix. We can do so, because we can assign each spatial location a location index from one to N , according to a certain rule, saying latitude from south to north, longitude from west to east, and elevation from low to high.

The following is a computing example of reducing a 3-dimensional array to a 2-dimensional array, i.e., a matrix.

```

828 #R code to reduce the dimension of an nD array
829 x <- array(1:(2*3*4), dim=c(2,3,4))
830 dim(x)
831 #[1] 2 3 4
832 x #a stack of four 2-by-3 matrices
833 #, , 1 #the first of the 3rd dim index
834
835 #[,1] [,2] [,3]
836 #[1,]    1     3     5
837 #[2,]    2     4     6
838 # ...
839 #install.packages('R.utils')
840 library(R.utils)
841 #flat all the other dim except the 3rd one
842 #flat the 1st and 2nd dim
843 y <- wrap(x, map=list(3, NA))
844 dim(y)
845 #[1] 4 6
846 y
847 #[,1] [,2] [,3] [,4] [,5] [,6]
848 #[1,]    1     2     3     4     5     6
849 #[2,]    7     8     9    10    11    12
850 #[3,]   13    14    15    16    17    18
851 #[4,]   19    20    21    22    23    24
852
853 #back to the original 3D array
854 array(t(y), dim = c(2,3,4))

```

855

1.4 A Set of Linear Equations

856

857 A somewhat different matrix example is the coefficient matrix of a system of linear equations. Solving a system of linear equations is very common in science and
 858 engineering. Finding numerical solutions for a climate model based on partial differential equations will usually involve solving large systems of linear equations.
 859

860 As an introduction to the coefficient matrix of linear equations, let us look at
 861 a simple elementary school mathematics problem: The sum of the ages of two
 862 brothers is 20 years and the difference of the ages is 4 years. What are the ages of
 863 the two brothers? One can easily guess that the older brother is 12 years old, and
 864 the younger one is 8.

865 If we form a set of equations, this would be

$$\begin{aligned} x_1 + x_2 &= 20 \\ x_1 - x_2 &= 4 \end{aligned} \tag{1.27}$$

866 where x_1 and x_2 stand for the brothers' ages.

867 The matrix form of these equations would be

$$Ax = b \tag{1.28}$$

869 which involves three matrices:

$$A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, b = \begin{bmatrix} 20 \\ 4 \end{bmatrix}. \quad (1.29)$$

870 Here, Ax means a matrix multiplication: $A_{2 \times 2}x_{2 \times 1}$.

871 The matrix notation of a system of two linear equations can be extended to
872 systems of many linear equations, hundreds or millions of equations in climate
873 modeling and climate data analysis. Typical linear algebra textbooks introduce
874 matrices in this way by describing linear equations in a matrix form. However, this
875 approach may be less intuitive for climate science, which emphasizes data. Thus,
876 our book uses data to introduce matrices as shown at the beginning of this chapter.

877 Although one can easily guess that the solution to the above simple matrix equation
878 (1.28) is $x_1 = 12$ and $x_2 = 8$, a more general method for computing the solution
879 may be using the R code shown below:

```
880 #Solve linear equations
881 A = matrix(c(1,1,1,-1), nrow = 2)
882 b = c(20,4)
883 solve(A, b)
884 #[1] 12 8 #This is the result x1=12, and x2=8.
```

885 This type of R program can solve more complicated systems of linear equations,
886 such as a system with 1,000 unknowns rather than two unknowns, as in this example.
887

888 The solution may be represented as

$$x = A^{-1}b, \quad (1.30)$$

889 where A^{-1} is the inverse matrix of A and was found earlier in eq. (1.26). One can
890 verify that

$$A^{-1}b = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} 20 \\ 4 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix} \quad (1.31)$$

891 is indeed the solution of the system of two linear equations.

892 1.5 Eigenvalues and eigenvectors of a square space 893 matrix

894 “Eigenvalue” is a partial translation of the German word “eigenwert,” meaning
895 “self-value” or “intrinsic value.” In German, “eigen” can mean “self” or “own”, as
896 in “one’s own,” and “wert” means “value.”

897 A square matrix is a matrix for which the number of rows is equal to the number
898 of columns. The matrix thus has the shape of a square. Similarly, other matrices
899 may be called rectangular matrices, or tall matrices.

901 1.5.1 Matrices of data anomalies, standardized anomalies, 902 covariance, and correlation 903

In climate science, one often considers the covariance or correlation among N stations, N grid boxes, or N grid points. The covariance between station i and station j is denoted by c_{ij} , and the corresponding correlation is denoted by r_{ij} . Then $C = [c_{ij}]_{N \times N}$ is the covariance matrix, and $Corr = [r_{ij}]_{N \times N}$ is the correlation matrix. They can be estimated by the observational or model data.

909 Let us use an $N \times Y$ matrix

$$X = [x_{it}]_{N \times Y} \quad (1.32)$$

to represent the climate data of these N stations with Y time steps. A time step can be a month, or a year, or a week, depending on the application needs and the data availability. The covariance and correlation matrices of the N stations for a climate parameter are square matrices. Both covariance and correlation matrices can be estimated using the anomaly data, which are defined as data minus their temporal mean of each station, i.e.,

$$A_{N \times Y} = [a_{it}]_{N \times Y} = [x_{it} - \bar{x}_i]_{N \times Y}, \quad (\text{Anomaly data matrix}), \quad (1.33)$$

916 where

$$\bar{x}_i = \frac{1}{Y} \sum_{t=1}^Y x_{it} \quad (1.34)$$

is the climate data mean of station i ($i = 1, 2, \dots, N$). The mean is also called climatology or climate normal of the climate parameter. Thus, the temporal mean of the anomaly data for each station should be zero. However, in many climate applications, the temporal data are not continuous and were missing in the earlier period or in the middle. It is impossible to compute the temporal mean for the entire period of length Y . Instead, the climatology is computed using a sub-interval of data period, whose length is smaller than Y , for example, using 30 years from 1971 to 2000. The 30-year climatology has almost become a recent community standard because of the period's best space-time coverage of observations. Consequently, the temporal mean of the anomaly data based on the 30-year climatology is often non-zero.

Another case is that some applications require us to compute climatologies using 10 years or shorter because of climate change. In any event, climatology and anomalies are commonly used terms in climate data analysis. Their exact definitions vary depending on practical applications. Chapter 11 of this book has more detailed discussions on climatology, anomalies, statistics, and visualization for the historical climate data which are incomplete with missing values.

When the anomaly matrix is normalized by the standard deviation of each station, then the result is called the standardized anomaly matrix:

$$A_{sd} = \left[\frac{x_{it} - \bar{x}_i}{s_i} \right]_{N \times Y}, \quad (\text{Standardized anomaly matrix}), \quad (1.35)$$

936 where

$$s_i = \left[\frac{1}{Y} \sum_{t=1}^Y (x_{it} - \bar{x}_i)^2 \right]^{1/2} \quad (1.36)$$

937 is the estimated standard deviation of station i ($i = 1, 2, \dots, N$). Strictly speaking
 938 from statistical theory, this standard deviation formula should use $Y - 1$ as the
 939 denominator, rather than simply Y , according to formula (B.4) in Chapter 3. How-
 940 ever, climate scientists often use Y , not $Y - 1$, perhaps because formula (1.36) has
 941 a clear meaning as the mean of square errors. In practice, the difference between
 942 the two formulas is small since Y is usually greater or equal to 30. Climate data
 943 analysis applications often use 30 years of data, say from 1971 to 2000, to compute
 944 climatology and standard deviation.

945 An advantage of using the standardized anomaly matrix is that it is dimension-
 946 less. A climate parameter, such as temperature, often has a larger variance over the
 947 land than ocean. The standardized anomaly data show better homogeneity between
 948 ocean and land.

949 Then, the covariance and correlation matrices can be estimated by the following
 950 formulas

$$C = \frac{1}{Y} A(A^t) \quad (\text{Covariance matrix}), \quad (1.37)$$

$$\text{Corr} = \frac{1}{Y} A_{sd}(A_{sd})^t \quad (\text{Correlation matrix}). \quad (1.38)$$

951 Once more, applications often use a 30-year mean, not the mean of the entire
 952 time period of length Y , to compute the covariance and correlation matrices.

953
 954 **Example 1.7** This example's data matrix A has $N = 2$ and $Y = 3$. The data's
 955 covariance matrix is `covm`.

```
956 #Spatial covariance matrix
957 dat = matrix(c(0,-1,1,2,3,4), nrow=3)
958 dat
959 colMeans(dat)
960 A = sweep(dat, 2, colMeans(dat))
961 A
962 # [,1] [,2]
963 #[1,] 0 -1
964 #[2,] -1 0
965 #[3,] 1 1
966 covm=(1/(dim(A)[2]))*A%*%t(A)
967 covm #is the covariance matrix.
968 # [,1] [,2] [,3]
969 #[1,] 0.5 0.0 -0.5
970 #[2,] 0.0 0.5 -0.5
971 #[3,] -0.5 -0.5 1.0
```

1.5.2 Eigenvectors and their corresponding eigenvalues

973 The covariance matrix times a vector u yields a new vector in a different direction.

```
976   u = c(1, 1, 0)
977   v = covm %*% u
978   v
979   #[,1]
980   #[1,] 0.5
981   #[2,] 0.5
982   #[3,] -1.0
983   #u and v are in different directions
```

984 In general, when we consider a given square matrix C and a given vector u , the
985 product Cu is usually not in the same direction as u , as shown in the above example.
986 However, there is always a “self-vector” vector w for each covariance matrix C such
987 that Cw is in the same direction as w , i.e., Cw and w are parallel, and this fact is
988 expressed as

$$Cw = \lambda w, \quad (1.39)$$

989 where λ is a scalar which has the property that it simply scales w so that the above
990 equation holds. This scalar λ is called an eigenvalue (i.e., a “self-value”, “own-
991 value”, or “characteristic value” of the matrix C), and w is called an eigenvector.

992 R can calculate the eigenvalues and eigenvectors of a covariance matrix `covm` with
993 a command `eigen(covm)`. The output is in an R data frame, which has two storages:
994 `ew$values` for eigenvalues and `ew$vector` for eigenvectors, as shown below.

```
995 #Eigenvectors of a covariance matrix
996 ew = eigen(covm)
997 ew
998 #$values
999 #[1] 1.500000e+00 5.000000e-01 1.332268e-15
1000
1001 #$vectors
1002 # [,1]      [,2]      [,3]
1003 #[1,] -0.4082483 -7.071068e-01 0.5773503
1004 #[2,] -0.4082483  7.071068e-01 0.5773503
1005 #[3,]  0.8164966  8.881784e-16 0.5773503
1006
1007 #Verify the eigenvectors and eigenvalues
1008 covm %*% ew$vectors[,1]/ew$values[1]
1009 # [,1]
1010 #[1,] -0.4082483
1011 #[2,] -0.4082483
1012 #[3,]  0.8164966
1013 #This is the first eigenvector
1014
1015 w = ew$vectors[,1] # is an eigenvector
```

1016 A 3×3 covariance matrix has three eigenvalues, and three eigenvectors (λ_1, w_1)
1017 (λ_2, w_2) , and (λ_3, w_3) . These are called eigenpairs.

1018 In this example,

$$\lambda_1 = 1.5, \quad w_1 = \begin{bmatrix} -0.4082483 \\ -0.4082483 \\ 0.8164966 \end{bmatrix}, \quad (1.40)$$

$$\lambda_2 = 0.5, \quad w_2 = \begin{bmatrix} -0.7071068 \\ 0.7071068 \\ 0 \end{bmatrix}, \quad (1.41)$$

$$\lambda_3 = 0.0, \quad w_3 = \begin{bmatrix} 0.5773503 \\ 0.5773503 \\ 0.5773503 \end{bmatrix}. \quad (1.42)$$

1019 The eigenvectors are frequently called modes, or empirical orthogonal functions
1020 (EOFs) in climate science. The term EOF was coined by Edward Lorenz (1917-
1021 2008) in his 1956 paper on statistical weather prediction (Lorenz 1956).

1022 Each eigenvalue is equal to the variance of the data projection on the corre-
1023 sponding eigenvector, and is thus positive. The sum of all the eigenvalues of such
1024 an $N \times N$ matrix represents the total variance of the climate system observed at
1025 these N stations. The first eigenvalue λ_1 is the largest, corresponding to the largest
1026 spatial variability of the climate field under study. The eigenvalues sizes follow the
1027 order $\lambda_1 \geq \lambda_2 \geq \dots$.

1028 However, in carrying out an analysis of climate data, one can often find the
1029 important patterns as eigenvectors more directly from the anomaly data matrix
1030 A without computing the covariance matrix C explicitly. This is known as the
1031 singular value decomposition (SVD) approach (Golub and Reinsch 1970), which we
1032 discuss next. It separates the space-time anomaly data into a space part, a time
1033 part, and what we may think of as a variation in energy part. This mathematical
1034 method of space-time decomposition is universally applicable to any data that we
1035 sample in space and time, and it can often help to develop physical insight and
1036 scientific understanding of the phenomena and their properties as contained in
1037 the observational data. Efficient computing methods of SVD have been extensively
1038 researched and developed since the 1960s. Gene H. Golub (1932-2007) was a leading
1039 researcher in this effort and is remembered as “Professor SVD” by his Stanford
1040 colleagues and the world mathematics community.

1041 1.6 An SVD representation model for space-time data

1042

1043 1.6.1 Space-Time Data Matrix and Its decomposition

1044

1045 We encounter space-time data every day, a simple example being the air tempera-
1046 ture at different locations at different times. If you take a plane to travel from San
1047 Diego to New York, you may experience the temperature at San Diego in the morn-

ing when you depart and that at New York in the evening after your arrival. Such data have many important applications. We may need to examine the precipitation conditions around the world at different days in order to monitor agricultural yields. A cellphone company may need to monitor its market share and the temporal variations of that quantity in different countries. A physician may need to monitor a patient's symptoms in different areas of the body at different times. The observed data in all these examples can form a space-time data matrix with the row position corresponding to the spatial location and the column position corresponding to time, as in Table 1.1.

Table 1.1 Space-time data table

| | Time 1 | Time 2 | Time 3 | Time 4 |
|---------|--------|--------|--------|--------|
| Space 1 | D11 | D12 | D13 | D14 |
| Space 2 | D21 | D22 | D23 | D24 |
| Space 3 | D31 | D32 | D33 | D34 |
| Space 4 | D41 | D42 | D43 | D44 |
| Space 5 | D51 | D52 | D53 | D54 |

Graphically, the space-time data may typically be plotted as a time series at each given spatial position, or as a spatial map at each given time. Although these straight-forward graphical representations can sometimes provide very useful information as input for signal detection, the signals are often buried in the data and may need to be detected by different linear combinations in space and time. Sometimes the data matrices are extremely large, with millions of data points in either space or time. Then the question arises as to how can we extract the essential information in such a big data matrix? Can we somehow manage to represent the data in a more simple and yet more useful way? A very useful approach to such a task involves a space-time separation. Singular value decomposition (SVD) is a method designed for this purpose. SVD decomposes a space-time data matrix into a spatial pattern matrix U , a diagonal energy level matrix D , and a temporal matrix V^t , i.e., the data matrix A is decomposed into

$$A_{N \times Y} = U_{N \times m} D_{m \times m} (V^t)_{m \times Y}. \quad (1.43)$$

where N is the spatial dimension, Y is the temporal length, $m = \min(N, Y)$, and V^t is the transpose of V . The columns of U are a set of spatial eigenvectors. They are orthonormal vectors, each of which is orthogonal to another and has its length equal to one. Thus, the word "orthonormal" comes from two words "orthogonal" and "normal." The orthonormal property can be expressed by the following formula:

$$\mathbf{u}_l \cdot \mathbf{u}_k = \delta_{lk}, \quad (1.44)$$

where \mathbf{u}_l and \mathbf{u}_k are column vectors of U , and δ_{lk} is the Kronecker delta equal to

1076 zero when $k \neq l$ and one when $k = l$. The columns of V are also a set of orthonormal
 1077 vectors known as temporal eigenvectors.

1078 Usually, the elements of the U and V matrices are unitless (i.e., dimensionless),
 1079 and the unit of the D elements is the same as the elements of the data matrix. For
 1080 example, if A is a space-time precipitation data matrix with a unit [mm/day], then
 1081 the dimension of the D elements is also [mm/day].

1082 **1.6.2 SVD of the Space-Time Anomaly Data Matrix and**
 1083 **Eigenvalue Problem of a Covariance Matrix**

1085
 1086 **Example 1.8** SVD for a 2×3 matrix.

```
1087 #Develop a 2-by-3 space-time data matrix for SVD
1088 A=matrix(c(1,-1,2,0,3,1), nrow=2)
1089 A
1090 #[,1] [,2] [,3]
1091 #[1,]    1    2    3
1092 #[2,]   -1    0    1
1093 #Perform SVD calculation
1094 msvd=svd(A)
1095 msvd
1096 msvd$d
1097 #[1] 3.784779 1.294390
1098 msvd$u
1099 #[,1]      [,2]
1100 #[1,] -0.9870875 -0.1601822
1101 #[2,] -0.1601822  0.9870875
1102 msvd$v
1103 #[,1]      [,2]
1104 #[1,] -0.2184817 -0.8863403
1105 #[2,] -0.5216090 -0.2475023
1106 #[3,] -0.8247362  0.3913356
1107 #One can verify that A=UDV', where V' is transpose of V.
1108 verim=msvd$u%*%diag(msvd$d)%*%t(msvd$v)
1109 verim
1110 #[,1]      [,2] [,3]
1111 #[1,] 1 2.000000e+00 3
1112 #[2,] -1 1.665335e-16 1
1113 round(verim)
1114 #[,1] [,2] [,3]
1115 #[1,] 1 2 3
1116 #[2,] -1 0 1
1117 #This is the original data matrix A
```

1118

1119 The covariance of the space-time anomaly matrix A is a spatial matrix:

$$C = \frac{1}{Y} AA^t, \quad (1.45)$$

1120 where Y is the number of columns of A and is equal to the length of time being
1121 considered.

1122 Strictly speaking, each column of an anomaly matrix A has a zero mean. But in
1123 practice, due to different normalization periods, the mean is not zero.

```
1124 covm = (1/(dim(A)[2]))*A%*%t(A)
1125 eigcov = eigen(covm)
1126 eigcov$values
1127 #[1] 4.7748518 0.5584816
1128 eigcov$vectors
1129 # [,1]      [,2]
1130 #[1,] -0.9870875 0.1601822
1131 #[2,] -0.1601822 -0.9870875
```

1132 Thus, the eigenvectors of a covariance matrix are the same as the SVD eigenvectors
1133 of the anomaly data matrix. The eigenvalues of the covariance matrix and the SVD
1134 have following relationship

```
1135 ((msvd$d)^2)/(dim(A)[2])
1136 #[1] 4.7748518 0.5584816
1137 eigcov$values
1138 #[1] 4.7748518 0.5584816
```

1139 Therefore, the EOFs from a given space-time dataset can be calculated directly by
1140 using an SVD command and do not need the step of calculating the covariance ma-
1141 trix. With efficient SVD algorithms, this shortcut can save significant amounts of
1142 time for an EOF analysis, also known as the principal component analysis (PCA) in
1143 the statistics community, compared to the traditional covariance matrix approach.
1144 Therefore, the result is extremely helpful for the EOF analysis, which is an indis-
1145 pensable modern tool of climate data analysis. The result is formally described as
1146 a theorem whose proof is also provided as follows.

1147 **Theorem 1.1** *The eigenvectors \mathbf{u}_k of the covariance matrix $C = (1/Y)AA^t$*

$$\mathbf{C}\mathbf{u}_k = \lambda_k \mathbf{u}_k, \quad (k = 1, 2, \dots, N), \quad (1.46)$$

1148 *are the same as the SVD spatial modes of $A = UDV^t$. The eigenvalues λ_k of $C_{N \times N}$
1149 and the SVD eigenvalues d_k of $A_{N \times Y}$ have the following relationship*

$$\lambda_k = d_k^2/Y \quad (k = 1, 2, \dots, Y), \quad \text{when } Y \leq N, \quad (1.47)$$

1150 *where Y is the total time length (i.e., time dimension) of the anomaly data matrix
1151 A , and N is the total number of stations for A (i.e., space dimension).*

1152 **Proof** The SVD of the space-time data matrix is

$$A = UDV^t. \quad (1.48)$$

1153 The data matrix A 's corresponding covariance matrix is thus

$$\begin{aligned} C &= \frac{1}{Y}AA^t \\ &= \frac{1}{Y}UDV^t(UDV^t)^t \\ &= \frac{1}{Y}UDV^t(VDU^t) \\ &= \frac{1}{Y}UD(V^tV)DU^t \\ &= \frac{1}{Y}UDIDU^t \\ &= \frac{1}{Y}UD^2U^t. \end{aligned} \tag{1.49}$$

1154 In the above, we have used $V^tV = I_Y$ is a $Y \times Y$ identity matrix according to
1155 the SVD definition. The identity matrix's dimension is Y if $Y \leq N$. Otherwise,
1156 $V^tV = I_N$.

1157 The expression $C = \frac{1}{Y}UD^2U^t$ is the EOF expansion of the covariance matrix
1158 and means that a covariance matrix consists of EOFs and its associated variance,
1159 or "energy."

1160 The covariance matrix's eigenvalue problem is

$$CU = \frac{1}{Y}UD^2U^tU = \frac{1}{Y}UD^2 = U\Lambda, \tag{1.50}$$

1161 where

$$\Lambda = \frac{1}{Y}D^2 \tag{1.51}$$

1162 is the diagonal eigenvalue matrix with its elements as

$$\lambda_k = d_k^2/Y \quad (k = 1, 2, \dots, Y). \tag{1.52}$$

1163 In eq. (1.50), we used $U^tU = I_Y$ based on the SVD definition. Equation (1.50)
1164 becomes an eigenvalues problem because Λ is a diagonal matrix: $C\mathbf{u}_k = \lambda_k\mathbf{u}_k$ ($k =$
1165 $1, 2, \dots, Y$), where \mathbf{u}_k is the k th column vector of the space matrix U . Thus, eq.
1166 (1.50) implies the first part of the theorem: The space eigenvectors U from the
1167 SVD of the space-time data matrix A are the same as the eigenvectors of the
1168 corresponding covariance matrix C .

1169 Equation (1.52) is exactly the second part of the theorem. The proof is thus
1170 complete.

1171 In the above proof, we implicitly assumed that the columns of data matrix A are
1172 independent when $Y \leq N$. Independence of a set of vectors means that no one can
1173 be expressed in terms of the others, so no vectors can be replaced. Real climate
1174 data always satisfy this independence assumption.

1175

1.7 Image analysis using SVD and R

1176

1177 This section shows an example of using SVD to analyze a color photo, The left
 1178 panel of Fig. 1.2 shows the color photo. The right panel of the figure is the black-
 1179 and-white photo produced by the grayscale data derived from the color photo. We
 1180 will also show how to reconstruct a photo from various SVD modes, as wells as how
 1181 to blend the monotone photos into a color photo.

1182
1183

1.7.1 Data for a color photo and its grayscale figure

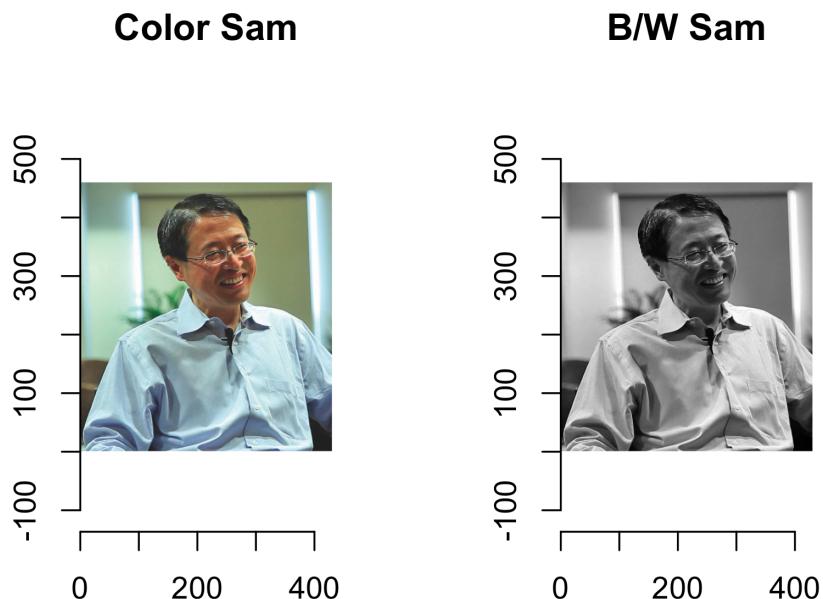


Fig. 1.2 Color and grayscale photos

1184 Figure 1.2 can be generated by the following computer code and data.

```
1185 #R code
1186 #Image analysis using SVD
1187 #Ref: imager: an R package for image processing
1188 #https://dahtah.github.io/imager/imager.html
1189
1190 setwd("~/mathmodel")
1191 #install.packages('imager') #run this if not installed yet
1192 library(imager)
```

```

1193 dat <- load.image('data/SamPhoto.png') #355 KB file size
1194 dim(dat)
1195 #[1] 430 460 1 3
1196 #430 rows and 460 columns, 430*460 = 197,800 pixels
1197 #1 photo frame, 3 RGB colors
1198 #If a video, 1 will become 150 frames or more
1199
1200 dat[1:3, 1:4,1,1]
1201 #Show part of the daat
1202
1203 #plot the color figure
1204 plot(dat,
1205   xlim = c(0, 430), ylim = c(0, 460),
1206   main = 'A Color Photo of Sam')
1207
1208 #Make the photo black-and-white
1209 graydat = grayscale(dat)
1210 dim(graydat)
1211 #[1] 430 460 1 1
1212 #430 rows and 460 columns, 1 photo frame, 1 grayscale [0, 1]
1213 #plot the gray b/w photo
1214 plot(graydat,
1215   xlim = c(0, 430), ylim = c(0, 460),
1216   main = 'B/W Sam')
1217 #Plot the color and b/w photos together
1218 par(mfrow = c(1, 2))
1219 plot(dat,
1220   xlim = c(0, 430), ylim = c(0, 460),
1221   main = 'Color Sam')
1222 plot(graydat,
1223   xlim = c(0, 430), ylim = c(0, 460),
1224   main = 'B/W Sam')
1225 dev.off()

```

1.7.2 SVD of the grayscale photo data

1.7.2.1 SVD singular values and their scree plot

1229 Data science and statistics often use the concept of how much variance is explained.
 1230 Hence, we plot the scree plot (see Fig. 1.3) based on variances, i.e., the eigenvalues
 1231 of a covariance matrix, which are the squares of SVD's singular values.

1232 The scree plot shows that the first mode is dominant, explaining around 94% of
 1233 variance. The variance contributions from other modes are almost zero. However,
 1234 these modes are extremely important in producing the details of the photo, as we
 1235 will see from the data reconstruction. This will be shown in the next subsection.

1236 Figure 1.3 and its associated SVD can be made using the following computer
 1237 code.

```

1238 #R code
1239 #SVD analysis of the grayscale data
1240 svdDat = svd(graydat)
1241 SVDD = svdDat$d
1242 percentD = 100*(SVDD^2)/sum(SVDD^2)

```

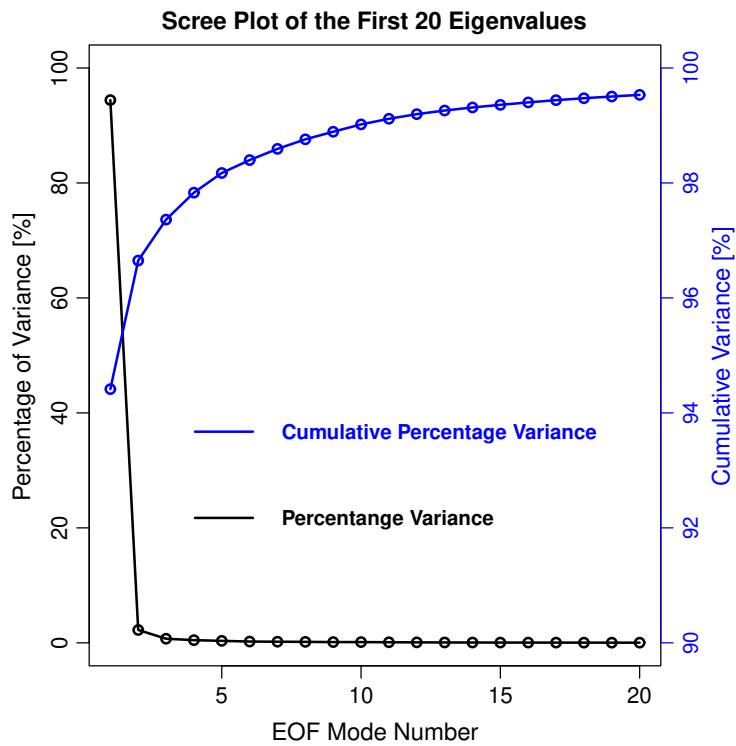


Fig. 1.3 Scree plot based on the grayscale photo's SVD singular values.

```

1243 cumpercentD = cumsum(percentD)
1244 modeK = 1:length(SVDD)
1245 dev.off()
1246 plot(modeK[1:20], percentD[1:20],
1247       type = 'o', col = 'blue',
1248       xlab = 'Mode_number', pch = 16,
1249       ylab = 'Percentage_of_mode_variance',
1250       main = 'Scree_Plot_of_SVD_B/W_Photo_Data')
1251 K = 20
1252 lam = (svdDat$d)^2
1253 lamK=lam[1:K]
1254 lamK
1255 #setEPS() #Plot the figure and save the file
1256 #postscript("fig0608.eps", width = 6, height = 4)
1257 par(mar=c(4,4,2,4), mgp=c(2.2,0.7,0))
1258 plot(1:K, 100*lamK/sum(lam), ylim=c(0,100), type="o",
1259       ylab="Percentage_of_Variance[%]",
1260       xlab="EOF_Mode_Number",
1261       cex.lab=1.2, cex.axis = 1.1, lwd=2,
1262       main="Scree_Plot_of_the_First_20_Eigenvalues")
1263 legend(3,30, col=c("black"),lty=1, lwd=2.0,
1264         legend=c("Percentange_Variance"), bty="n",
1265         text.font=2,cex=1.0, text.col="black")
1266 par(new=TRUE)

```

```

1267 plot(1:K, cumsum(100*lamK/sum(lam)),
1268   ylim = c(90,100), type="o",
1269   col="blue", lwd=2, axes=FALSE,
1270   xlab="", ylab="")
1271 legend(3,80, col=c("blue"), lty=1, lwd=2.0,
1272   legend=c("Cumulative\u2022Percentage\u2022Variance"), bty="n",
1273   text.font=2, cex=1.0, text.col="blue")
1274 axis(4, col="blue", col.axis="blue", mgp=c(3,0.7,0))
1275 mtext("Cumulative\u2022Variance\u2022[%]", col="blue",
1276   cex=1.2, side=4, line=2)

```

1.7.2.2 The grayscale photo reconstruction from the SVD modes

Figure 1.4 shows four reconstructions from different SVD modes. When using all the modes, the reconstruction is 100% accurate and completely recovers the original grayscale photo. When using the first 20 modes, the reconstructed photo is fuzzy. Although 99% of variance is explained in the first 20 modes, many details are still contained in the higher modes whose variance contribution is less 1%. When using the three modes, the reconstructed photo is so fuzzy that it is barely recognizable until you look at it carefully in the direct front view. The first three modes explain almost 98% of total variance according to the scree plot. Many fields of natural sciences use only the modes for 90% variance. This practice may need careful scrutiny when more detailed signals of small scales, such as this photo analysis, are important. When using the 80 modes between 21st and the 100th modes, the reconstruction shows the person with glasses without the background layout. These 80 modes account for less than 1% variance but can provide this much information. This seems surprising.

Using the computer code we provide here, you can easily test reconstructions using different modes. You can also use your preferred photos and other figures. You may come out with some surprising discoveries.

Figure 1.4 can be generated by the following computer code.

```

1296 #R code
1297 dev.off()
1298 par(mfrow = c(2, 2))
1299 kR = 430 #Recon from all 413 modes
1300 R430 = as.cimg(U[,1:kR] %*% D[1:kR, 1:kR] %*% t(V[, 1:kR]))
1301 plot(R430, main = "All\u2022430\u2022modes")
1302 kB = 20 #The first 20 modes
1303 R20 = as.cimg(U[,1:kB] %*% D[1:kB, 1:kB] %*% t(V[, 1:kB]))
1304 plot(R20, main = "The\u2022first\u202220\u2022modes")
1305 k = 3 #Recon from the first 3 modes
1306 R3 = as.cimg(U[,1:k] %*% D[1:k, 1:k] %*% t(V[, 1:k]))
1307 plot(R3, main = "The\u2022first\u20223\u2022modes")
1308 k1 = 21; k2 = 100 #Recon from 21st to 100th modes
1309 Rk1_k2 = as.cimg(U[,k1:k2] %*% D[k1:k2, k1:k2] %*% t(V[, k1:k2]))
1310 plot(Rk1_k2, main = "21st\u2022to\u2022100th\u2022modes")
1311 dev.off()

```

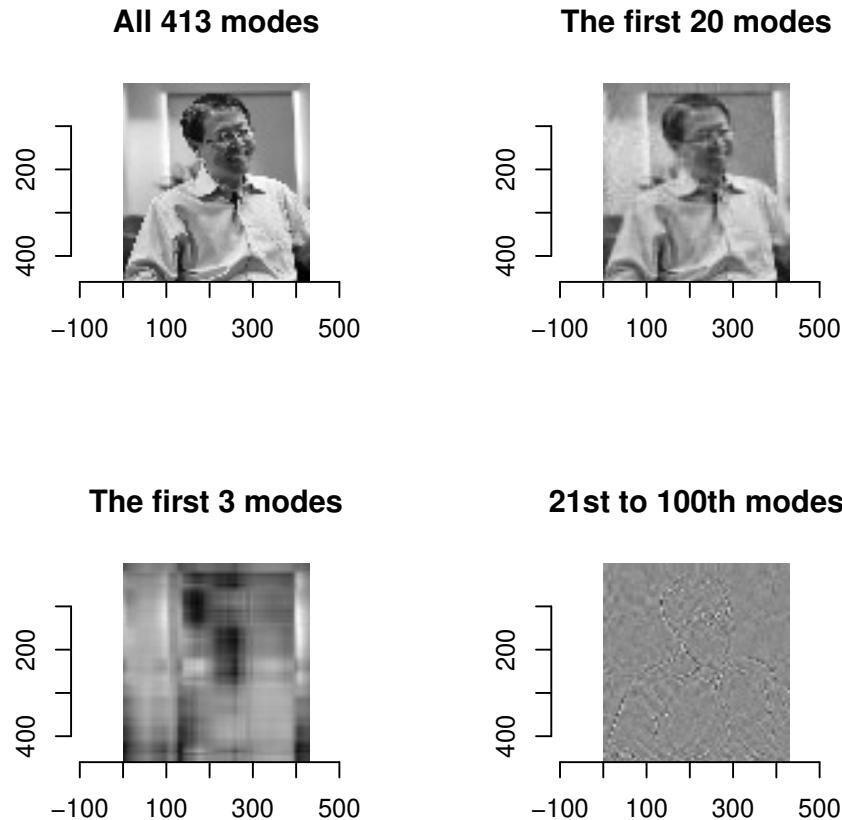


Fig. 1.4 Reconstruction of the grayscale photo using different numbers of SVD modes.

1312 **1.7.3 Reconstructing a color photo from the data of monotone**
 1313 **photos**
 1314

1315 We described earlier that the data for a color photo is a 4D array. The first two
 1316 dimensions are defined on the x-y pixels. The third dimension is the number of
 1317 photos. The fourth dimension has 3 entries for three monotone colors, such as R
 1318 (red), G (green), and B (blue). This subsection shows that you can recover the
 1319 original photo from the data for each of the three monotone figures. Figure 1.5
 1320 displays the recovered color photo and the three monotone photos.

1321 Figure 1.5 can be generated by the following computer code.

```
1322 #Rcode
1323 #Three monotone photos and their color photo
1324 dim(dat) #4D array of a color photo
```

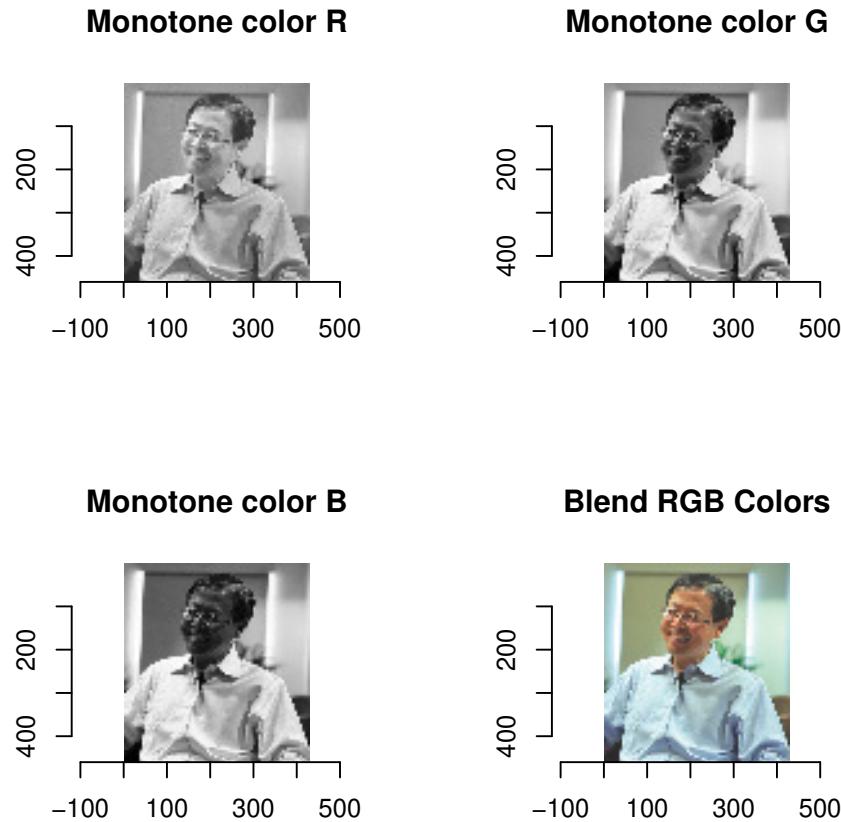


Fig. 1.5 Three monotone photos and the color photo recovered from them.

```

1325 # [1] 430 460    1    3
1326 dev.off()
1327 par(mfrow = c(2, 2))
1328 R = as.cimg(apply(t(dat[, 1, 1]), 1, rev))
1329 plot(R, main = 'Monotone_color_R')
1330 G = as.cimg(apply(t(dat[, 1, 2]), 1, rev))
1331 plot(G, main = 'Monotone_color_G')
1332 B = as.cimg(apply(t(dat[, 1, 3]), 1, rev))
1333 plot(B, , main = 'Monotone_color_B')
1334 #Bind the three monotone photos into one
1335 trippy = imappend(list(R,G,B), "c")
1336 dim(trippy) #color figure data
1337 #[1] 430 460    1    3
1338 plot(trippy, main ='Blend_RGB_Colors')
1339 dev.off()

```

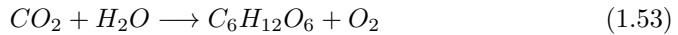
1340 This example shows a way to color a grayscale photo by blending additional

1341 monotone photos. Then, the difficulty is where to color and what color. AI coloring
1342 for old black-and-white films has been made in the market. AI coloring technology
1343 still needs improvement. Can SVD help?

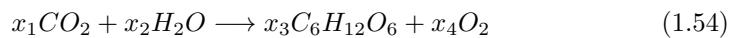
1344 1.8 Mass balance for chemical equations in marine 1345 chemistry

1346 This section describes another application of linear algebra in climate science. Ma-
1347 rine chemistry or atmospheric chemistry involves various kinds of chemical reaction
1348 equations which describe the conservation of mass during the chemical reactions.
1349 Then, the problem is how to systematically achieve the mass balance, i.e., to deter-
1350 mine the numbers of molecules on each side of an equation that depicts a chemical
1351 reaction. Here we use photosynthesis as an example to illustrate a linear algebra
1352 approach to this problem.

1353 In the process of photosynthesis, plants convert the solar radiant energy carried by
1354 photons, plus carbon dioxide (CO_2) and water (H_2O), into glucose ($\text{C}_6\text{H}_{12}\text{O}_6$) and
1355 oxygen (O_2). The chemical equation for this conversion could be written schematic-
1356 ally as



1357 However, the conservation of mass requires that the atomic weights on both sides
1358 of the equation be equal. The photons have no mass. Thus, the chemical equation
1359 as written above is incorrect. The correct equation should specify precisely how
1360 many CO_2 molecules react with how many H_2O molecules to generate how many
1361 $\text{C}_6\text{H}_{12}\text{O}_6$ and O_2 molecules. Suppose that these coefficients are x_1, x_2, x_3, x_4 . We
1362 then have



1363 Making the number of atoms of carbon on the left and right sides of the equation
1364 equal yields

$$x_1 = 6x_3 \quad (1.55)$$

1365 because water and oxygen contain no carbon. Doing the same for hydrogen atoms
1366 leads to

$$2x_2 = 12x_3. \quad (1.56)$$

1367 Similarly, the balance of oxygen atoms is

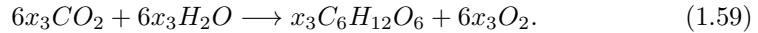
$$2x_1 + x_2 = 6x_3 + 2x_4. \quad (1.57)$$

1368 We thus have three equations in four variables. Thus, these equations have infinitely
1369 many solutions. We can set any variable fixed, and express the other three using

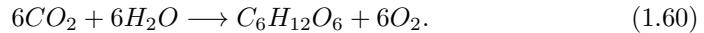
¹³⁷¹ this fixed variable. Since the largest molecule is glucose, we set its coefficient x_3
¹³⁷² fixed. Then we have

$$x_1 = 6x_3, \quad x_2 = 6x_3, \quad x_4 = 6x_3. \quad (1.58)$$

¹³⁷³ Thus, the chemical equation is



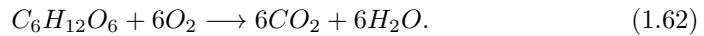
¹³⁷⁴ If we want to produce one glucose molecule, i.e., $x_3 = 1$, then we need 6 carbon
¹³⁷⁵ dioxide and 6 water molecules:



¹³⁷⁶ Similarly, one can write chemical equations for many common reactions, such as
¹³⁷⁷ iron oxidation



¹³⁷⁸ and the redox reaction in a human body which consumes glucose and converts the
¹³⁷⁹ glucose into energy, water and carbon dioxide



¹³⁸⁰ 1.9 Multivariate linear regression using matrix ¹³⁸¹ notations

¹³⁸²

¹³⁸³ This section is an application of matrix algebra in statistical data analysis, par-
¹³⁸⁴ ticularly on a linear regression for more than one variable. The linear regression
¹³⁸⁵ section discussed the fitting of a straight line on the xy -plane $y = b_1x + b_0$ to a pair
¹³⁸⁶ of data vectors of length N : $[x_d]_{N \times 1}, [y_d]_{N \times 1}$. It resulted in correlation, trend and
¹³⁸⁷ other regression quantities. An example is the correlation between the January SOI
¹³⁸⁸ as x and U.S. temperature as y . The non-trivial correlation can suggest that there
¹³⁸⁹ may be a physical mechanism to explain how the January SOI influences the U.S.
¹³⁹⁰ January temperature.

¹³⁹¹ Here, we use $y = b_1x + b_0$ as the representation of a deterministic fitting func-
¹³⁹² tion. This expression does not involve random variables. Most statistics books would
¹³⁹³ consider random linear models, which distinguishes between random variables and
¹³⁹⁴ their deterministic estimators by data. Our simple linear mathematical model for-
¹³⁹⁵ mulation here is equivalent to using only the deterministic estimators.

¹³⁹⁶ The U.S. January temperature can be influenced by multiple factors in addition
¹³⁹⁷ to the SOI. These factors may include the North Atlantic sea surface temperature
¹³⁹⁸ (SST), the North Pacific SST, Arctic pressure conditions, etc. Then, the linear model
¹³⁹⁹ becomes

$$y = b_0 + b_1x_1 + b_2x_2 + \cdots + b_nx_n. \quad (1.63)$$

¹⁴⁰⁰ When $n = 1$, this degenerates into the single variable regression.

1401 In the following we will present a few R examples of multivariate regression and
 1402 its applications. The mathematical theory behind the R code for a multivariate
 1403 regression deals mostly with the matrix operations, which can be found from any
 1404 standard textbooks covering multivariate regression.

1405

1406 **Example 1.9** This example shows a two-variable regression.

$$y = b_0 + b_1 x_1 + b_2 x_2. \quad (1.64)$$

1407 Geometrically, this is an equation of a plane in a 3D space (x_1, x_2, y) . Given three
 1408 points not on a straight line, a plane can be determined uniquely. This means spec-
 1409 ifying three x_1 coordinate values, three x_2 coordinate values and three y coordinate
 1410 values, which can be done by means of the following R code:

```
1411 x1=c(1,2,3) #Given the coordinates of the 3 points
1412 x2=c(2,1,3)
1413 y=c(-1,2,1)
1414 df=data.frame(x1,x2,y) #Put data into the data.frame format
1415 fit <- lm(y ~ x1 + x2, data=df)
1416 fit#Show the regression results
1417 #Call:
1418 #  lm(formula = y ~ x1 + x2, data = df)
1419 #Coefficients:
1420 # (Intercept)      x1          x2
1421 #-5.128e-16   1.667e+00  -1.333e+00
1422
1423 1.667*x1-1.333*x2 #Verify that 3 points determining a plane
1424 #[1] -0.999 2.001 1.002
```

1425 **Example 1.10** This example will show that four arbitrarily specified points can-
 1426 not all be on a plane. The fitted plane has the shortest distance squares, i.e., the
 1427 least squares (LS), or minimal mean square error (MMSE). Thus, the residuals are
 1428 non-zero, in contrast to the zero residuals in the previous example.

```
1429 u=c(1,2,3,1)
1430 v=c(2,4,3,-1)
1431 w=c(1,-2,3,4)
1432 mydata=data.frame(u,v,w)
1433 myfit <- lm(w ~ u + v, data=mydata)
1434 summary(myfit)#Show the result
1435 #Coefficients:
1436 # Estimate Std. Error t value Pr(>|t|)
1437 #(Intercept) 1.0000    1.8708   0.535    0.687
1438 #u           2.0000    1.2472   1.604    0.355
1439 #v          -1.5000    0.5528  -2.714    0.225
```

1440 **Example 1.11** This example will show a general multivariate linear regression
 1441 using R. It has three independent variables, one dependent variable, and ten data
 1442 points. For R program simplicity, the data are generated by an R random number
 1443 generator. Again, R requires that the data be put into data frame format so that

1444 a user can clearly specify which are independent variables, also called explainable
 1445 variables, and which is the dependent variable, also called response variable.

```
1446 dat=matrix(rnorm(40),nrow=10, dimnames=list(c(letters[1:10]), c(LETTERS[23:26])))
1447 fdat=data.frame(dat)
1448 fit=lm(Z~ W + X + Y, data=fdat)
1449 summary(fit)
1450
1451 #Coefficients\index{linear regression!coefficients}
1452 #             Estimate Std. Error t value Pr(>|t|)
1453 #(Intercept) 0.36680   0.16529   2.219   0.0683
1454 #W           0.11977   0.20782   0.576   0.5853
1455 #X           -0.53277  0.19378  -2.749   0.0333
1456 #Y           -0.04389  0.14601  -0.301   0.7739
```

1457 Thus, the linear model is

$$Z = 0.37 + 0.12W - 0.53X - 0.04Y. \quad (1.65)$$

1458 The 95% confidence interval for W 's coefficient is $0.12 \pm 2 \times 0.21$, that for X 's
 1459 coefficient is $-0.53 \pm 2 \times 0.19$, Y 's coefficient is $-0.04389 \pm 2 \times 0.15$. Each confi-
 1460 dence interval includes zero. Thus, there is no significant non-zero trend for the Z
 1461 data with respect to W, X, Y . This result is to be expected, because the data are
 1462 randomly generated and thus should not have a trend.

1463 In practical applications, a user can simply convert the data into the same data
 1464 frame format as shown here. Then, R command
 1465 `lm(formula = Z ~ W + X + Y, data = fdat)`
 1466 can do the regression job.

1468 R can also do nonlinear regression with specified functions, such as quadratic
 1469 functions and exponential functions. See examples from the URLs

1470 <https://www.zoology.ubc.ca/~schluter/R/fit-model/>
 1471 <https://stat.ethz.ch/R-manual/R-devel/library/stats/html/nls.html>

1.10 Chapter summary

1472 This chapter introduces matrices and linear algebra from the perspective of space-
 1473 time climate data, with rows corresponding to spatial locations (e.g., the order of
 1474 grid boxes, grid point IDs, or climate station IDs) and columns corresponding to
 1475 time (e.g., the month stamps January 1901, February 1901, ..., December 2017
 1476 for the monthly data, or the day stamps 1 June 2017, 2 June 2017, ..., 30 June
 1477 December 2017 for the daily data of June 2017). This way of introducing matrices
 1478 as organized collections of data is different from that of most textbooks which
 1479 describe a matrix as originating from a set of linear equations. This particular way
 1480 to arrange a space-time climate data matrix follows the similar philosophy of the
 1481

1483 netCDF data file, which has recently become a very popular format for representing
 1484 extensive climate datasets.

1485 The most essential methods of linear algebra have to do with matrix properties
 1486 (e.g., eigenvalues, eigenvectors, and SVD matrices) and matrix operations (e.g.,
 1487 solving linear equations, and making linear transformations). Our chapter has de-
 1488 scribed the following linear algebra methods:

- 1489 (i) Matrix arithmetic: addition of a matrix plus another matrix, subtraction of
 1490 a matrix minus another matrix, multiplication of a matrix times another
 1491 matrix, and multiplication of a matrix by a scalar, and division of the
 1492 identity matrix by a matrix which is a way to define the inverse matrix of
 1493 the latter.
- 1494 (ii) An eigenvector w of a matrix C is such a special vector that Cw is in w 's
 1495 own direction. Thus, there exists a scalar λ such that $Cw = \lambda w$, where λ
 1496 is called an eigenvalue of the matrix C corresponding to the eigenvector w .
- 1497 (iii) SVD decomposes the space-time data matrix A into three matrices: a dimen-
 1498 sionless spatial orthogonal matrix U as spatial eigenvectors, a dimension-
 1499 less temporal orthogonal matrix V as temporal eigenvectors, and a diagonal
 1500 “energy level” matrix D related to the eigenvalues of A 's covariance matrix,
 1501 i.e.,

$$A = UDV^t. \quad (1.66)$$

1502 The column vectors of U may represent the intrinsic spatial patterns of
 1503 the climate data (e.g., the warming pattern of the eastern tropical Pacific
 1504 for the El Niño based on the SST data), and those of V may represent
 1505 the intrinsic temporal patterns (e.g., the El Niño peaks in the time series
 1506 determined by a column vector of V). See the specific patterns shown in
 1507 the examples of the section on the SVD analysis for the reanalysis data.
 1508 The diagonal values of D show the “energy level” (or the strength) of the
 1509 corresponding spatial and temporal patterns. In the last few years, the SVD
 1510 method has become more widely used in science and engineering because of
 1511 its universal power in decomposing any space-time matrix. Mathematicians
 1512 have now begun to modernize the proofs of the theorems of linear algebra
 1513 by using the SVD approach, instead of the traditional approach known as
 1514 echelon reduction through row operations.

- 1515 (iv) Mathematical models based on linear equations have been developed and
 1516 employed for many climate science applications, such as the mass balance
 1517 model for chemical equations.
- 1518 (v) Multivariate linear regression has been formulated here from the perspective
 1519 of matrix algebra.
- 1520 (vi) R codes have been written for all these methods, which can be conveniently
 1521 used for solving climate data analysis problems in research and practical
 1522 applications.

1523 Although the basic materials of linear algebra appeared as early as the 17th

1524 century, widespread university classroom education in linear algebra in the United
1525 States began only in the 1950s and later (Tucker 1993). At that time, electronic
1526 computers began to become available to solve large numbers of linear equations for
1527 engineering applications. Today's vastly increased computing power now provides
1528 opportunities to modernize educational aspects of this branch of mathematics, such
1529 as using SVD to compute various properties of a matrix and to prove numerous
1530 theorems of linear algebra. Climate mathematics takes full advantage of this tech-
1531 nical progress. Many graduate students in climate science now use EOFs and PCs
1532 in their research, whereas the EOF-PC techniques were rarely taught as part of
1533 climate science education before the 1990s.

References and Further Readings

- 1535 [1] Golub, G.H. and C. Reinsch, 1970: Singular value decomposition and least
 1536 squares solutions. *Numerische mathematik* 14, 403-420.

This seminal paper established an important method, known as the Golub-Reinsch algorithm, to compute the eigenvalues of a covariance matrix from a space-time matrix A without actually first computing the covariance matrix AA^t . This algorithm makes the SVD computation very efficient, which helps scientists consider SVD as a genuine linear algebra method, not a traditionally regarded statistical method based on a covariance matrix.

- 1537
 1538 [2] Lorenz, E.N., 1956: Empirical orthogonal functions and statistical weather
 1539 prediction. *Scientific Report No. 1, Statistical Forecasting Project*. Air Force
 1540 Research Laboratories, Office of Aerospace Research, USAF, Bedford, Mas-
 1541 sachusetts, 49pp.

Edward N. Lorenz (1917-2008) was an American meteorologist and applied mathematician, who has been best known for his theories of chaos and the butterfly effect. Although he was not the first to invent the mathematical method of empirical orthogonal functions (EOFs), Lorenz, together with other scientists, such as Gerald R. North, made EOF a very popular method for climate data analysis.

- 1542
 1543 [3] Strang, G., 2016: *Introduction to Linear Algebra*. 5th edition, Wellesley-
 1544 Cambridge Press, Wellesley, MA 02482, 574pp.

Gilbert Strang (1934-) is an American mathematician and educator. His textbooks and pedagogy have been internationally influential. This text is one of the very few basic linear algebra books that includes excellent materials on SVD, probability, and statistics.

- 1545
 1546 [4] Tucker, A., 1993: The growing importance of linear algebra in undergraduate
 1547 mathematics. *College Mathematics Journal* 24, 3-9.

This paper describes the historical development of linear algebra, such as the term “matrix” being coined by J.J. Sylvester in 1848, and pointed out that “tools of linear algebra find use in almost all academic fields and throughout modern society.” The use of linear algebra in the big data era is now even more popular.

1548

Exercises

1549

1550

1551 **1.1** The following are the SVD results

```

1552     mat
1553         [,1] [,2]
1554 [1,]    1    1
1555 [2,]    1   -1
1556
1557     svd(mat)
1558 $d
1559 [1] 1.414214 1.414214
1560
1561 $u
1562             [,1]      [,2]
1563 [1,] -0.7071068 -0.7071068
1564 [2,] -0.7071068  0.7071068
1565
1566 $v
1567         [,1] [,2]
1568 [1,]    -1    0
1569 [2,]     0   -1

```

1570 Use $A = UDV^t$ to recover the first column of

```

1571     mat
1572         [,1] [,2]
1573 [1,]    1    1
1574 [2,]    1   -1

```

1575 Show detailed calculations of all the relevant matrices and vectors, and use
1576 space-time decomposition to describe your results.

1577 For extra credit: Describe the spatial and temporal modes, and their corre-
1578 sponding variances or energies.

- 1579 **1.2** Use R and the updated Darwin and Tahiti standardized SLP data to repro-
 1580 duce the EOFs and PCs and to plot the EOF pattern maps and PC time
 1581 series.
- 1582 **1.3** Do the same procedures in the previous problem but for original non-standardized
 1583 data. Comment on the difference of the results from the previous problem.
- 1584 **1.4** (a) Download the monthly precipitation data at five different stations over
 1585 the United States from the USHCN website:
- 1586
- 1587 http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn_map_interface.html
- 1588
- 1589 (b) Use R to organize the January data from 1951 to 2010 into the space-time
 1590 format.
- 1591 (c) Compute the climatology of each station as the 1971-2000 mean.
- 1592 (d) Compute the space-time anomaly data matrix A as the original space-time
 1593 data matrix minus the climatology.
- 1594 (e) Use R to make the SVD decomposition of the space-time anomaly data
 1595 matrix $A = UDV^t$.
- 1596 (f) Output the U and D matrices.
- 1597 **1.5** In the previous problem, use R and the formula UDV^t to reconstruct the
 1598 original data matrix A . This is a verification of the SVD decomposition, and
 1599 is also called EOF-PC reconstruction.
- 1600 **1.6** Use R to plot the maps of the first three EOF modes from the U matrix in
 1601 the previous problem in a way similar to the two El Niño mode maps you can
 1602 find on the Internet. Try to explain the meaning of climate in the EOF maps.
- 1603 **1.7** Use R to plot the first three PC time series from the V matrix in Problems
 1604 4.4 and 4.5. Try to explain the climate meaning of the time series.
- 1605 **1.8** (a) A covariance matrix C can be computed from a space-time observed
 1606 anomaly data matrix X which has N rows for spatial locations and Y columns
 1607 for time in years:

$$C = X \cdot X^t / Y \quad (1.67)$$

1608 This is an $N \times N$ matrix. Choose an X data matrix from the USHCN
 1609 annual total precipitation data at three California stations from north to
 1610 south [Berkeley, CA (040693); Santa Barbara, CA (047902); Cuyamaca, CA
 1611 (042239)] and five years from 2001 to 2005. Consider the anomaly data with
 1612 respect to the 2001-2005 mean and use R to calculate a covariance matrix for
 1613 $N = 3$ and $Y = 5$.

1614 (b) Use R to find the inverse matrix of the covariance matrix C .

1615 (c) Use R to find the eigenvalues and eigenvectors of C .

1616 (d) Use R to make SVD decomposition of the data matrix $X = UDV^t$. Ex-
 1617 plicitly write out the three matrices U , D and V .

1618 (e) Use R to explore the relationship between the eigenvalues of C and the
 1619 matrix D .

- 1620 (f) Compare the eigenvectors of C and the matrix U .
1621 (g) Plot the PC time series and describe their behavior.
1622
1623 **1.9** The burning of methane (CH_4) with oxygen (O_2) produces water (H_2O) and
1624 carbon dioxide (CO_2). Balance the chemical reaction equation.
1625 **1.10** The burning of propane (C_3H_8) with oxygen (O_2) produces water (H_2O) and
1626 carbon dioxide (CO_2). Balance the chemical reaction equation.
1627 **1.11** The burning of gasoline (C_8H_{18}) with oxygen (O_2) produces water (H_2O)
1628 and carbon dioxide (CO_2). Balance the chemical reaction equation.

2

Matrix Theory and Visualization

1631

1632 This chapter includes a slightly more advanced theory of matrix compared with the
 1633 last. It includes (i) the concepts of independence, (ii) spanned spaces by multiple
 1634 vectors, (iii) rank and other properties of a matrix, (iv) more on SVD, and (v)
 1635 visualization of matrices and their decomposed vectors and singular values.

1636 Many mathematicians regard mathematics as “beautiful,” using terms such as
 1637 “elegant,” “deep,” and “general.” Our book, however, pays more attention to *relevant*,
 1638 *useful*, and *modern* (*RUM*) from the perspectives of both mathematical sci-
 1639 ences and other fields. *Relevant* is reflected by that every matrix and its operations
 1640 can have an interpretation story that can be easily understood by a layman. *Useful*
 1641 is reflected by that each theory or method has non-trivial application examples in
 1642 science or engineering. *Useful* is reflected in the extensive use of matrices in modern
 1643 data science, machine learning, R or Python programming, which were not the case
 1644 two or three decades ago. Matrix visualization is an example of modernness.

1645 This chapter continues to feature the space-time data arrangement, which uses
 1646 rows of a matrix for spatial locations, and columns for temporal steps. This is the
 1647 universal and fundamental information structure of our world. The singular value
 1648 decomposition (SVD) helps reveal the spatial and temporal features of climate
 1649 dynamics as singular vectors and the strength of their variability as singular values.

1650

1651

2.1 Matrix Definitions

1652 A matrix is a rectangular array of numbers (or even expressions), often denoted by
 1653 an upper case letter in either boldface or plain **A** or *A*:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ a_{n1} & a_{n2} & \cdots & a_{np} \end{bmatrix} \quad (2.1)$$

1654 This is an $n \times p$ matrix, in which a_{ij} are called elements or entries of the matrix **A**, i
 1655 is the row index from 1 to n and j is the column index from 1 to p . The dimensions
 1656 of a matrix are denoted by subscript, e.g., $\mathbf{A}_{n \times p}$ indicating that the matrix **A** has n
 1657 rows and p columns. The matrix may also be indicated by square brackets around

1658 a typical element $\mathbf{A} = [a_{ij}]$ or sometimes $\{A\}_{ij}$, maybe even A_{ij} . If $n = p$, then the
 1659 the array is a square matrix.

1660 Figure 2.1 is an example of a space-time climate data matrix. It is a subset of the
 1661 $5^\circ \times 5^\circ$ gridded monthly surface air temperature anomalies from the NOAA Merged
 1662 Land Ocean Global Surface Temperature Analysis (NOAAGlobalTemp) (Version
 1663 4.0). The rows are indexed according to the spatial locations prescribed by the
 1664 latitude and longitude of the centroid of a $5^\circ \times 5^\circ$ grid box (See the entries of the
 1665 first two columns in boldface). The columns are indexed according to time (See
 1666 the first row entries in boldface). The other entries are the temperature anomalies
 1667 with respect to the 1971-2000 monthly climatology. The anomalies are arranged
 1668 according to the locations by rows and the time by columns. The units for the
 1669 anomaly data are $^\circ\text{C}$.

1670 For a given month, the spatial temperature data on the Earth's surface is itself a
 1671 2-dimensional array. To make a space-time matrix, we assign each grid box a unique
 1672 index s from 1 to n if the spatial region has n grid boxes. The index assignment is
 1673 subjective, depending on the application needs. The commonly used way is to fix a
 1674 latitude and increase the index number as the longitude increases, as indicated by
 1675 the first two columns of the data matrix shown in Figure 2.1. When the longitude
 1676 is finished at this latitude band, go to the next latitude band until the completion
 1677 of the latitude. This can go from south to north, from north to south. Of course,
 1678 one can fix the longitude first, and increase the index according to the ascending
 1679 or descending order of latitudes.

| Lat | Lon | 1934-3 | 1934-4 | 1934-5 | 1934-6 | 1934-7 | 1934-8 | 1934-9 | 1934-10 | 1934-11 | 1934-12 | 1935-1 | 1935-2 | 1935-3 | 1935-4 | 1935-5 |
|------|--------------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|--------|--------|--------|--------|--------|
| 32.5 | 242.5 | 1.86 | 1.14 | 1.03 | -0.65 | 0.12 | -0.27 | -0.30 | -0.30 | 0.13 | 0.34 | -0.48 | -0.18 | -1.43 | -0.50 | -0.84 |
| 32.5 | 247.5 | 3.14 | 2.42 | 2.29 | -2.08 | 0.93 | -0.09 | -0.49 | 0.46 | 0.21 | 0.79 | 0.07 | -0.18 | -2.04 | -0.44 | -2.32 |
| 32.5 | 252.5 | 1.42 | 1.94 | 2.47 | -0.24 | 1.18 | 1.04 | 0.11 | 1.30 | 0.46 | 0.73 | 0.93 | -1.07 | -0.09 | 0.18 | -2.32 |
| 32.5 | 257.5 | -0.97 | 0.96 | 0.95 | 1.89 | 1.44 | 1.49 | 0.50 | 2.68 | 1.55 | 0.87 | 2.65 | -0.36 | 1.95 | 0.64 | -2.18 |
| 32.5 | 262.5 | -1.89 | 1.09 | 0.51 | 2.64 | 2.12 | 2.36 | 0.07 | 2.67 | 1.90 | 0.47 | 2.40 | 0.14 | 2.50 | 0.04 | -1.66 |
| 32.5 | 267.5 | -1.36 | 0.82 | -0.06 | 1.36 | 0.89 | 1.46 | -0.55 | 2.28 | 1.19 | -0.17 | 2.14 | 0.54 | 3.27 | 0.21 | -0.35 |
| 32.5 | 272.5 | -0.98 | 0.61 | 0.26 | 0.82 | 0.36 | 0.38 | -0.15 | 1.37 | 0.93 | -0.66 | 1.31 | 0.23 | 2.45 | 0.76 | 0.85 |
| 32.5 | 277.5 | -1.26 | 0.51 | -0.24 | 0.75 | 0.65 | 0.39 | 0.71 | 0.91 | 0.36 | -0.85 | 0.99 | -0.25 | 2.12 | 0.68 | 0.90 |
| 32.5 | 282.5 | 0.54 | 0.88 | 0.09 | 0.25 | 0.32 | 0.13 | 0.20 | 1.08 | 0.51 | 0.55 | 0.81 | 0.87 | 1.31 | 0.94 | 0.77 |
| 32.5 | 287.5 | 0.72 | 0.99 | 0.29 | 0.39 | 0.36 | 0.12 | 0.50 | 1.04 | 0.32 | 0.38 | 0.23 | 0.46 | 0.62 | 0.58 | 0.39 |
| 32.5 | 292.5 | 0.79 | 0.93 | 0.27 | 0.48 | 0.23 | 0.18 | 0.90 | 1.01 | 0.48 | 0.22 | -0.23 | 0.11 | 0.21 | 0.30 | -0.04 |
| 32.5 | 297.5 | 0.68 | 0.59 | 0.26 | 0.33 | 0.17 | 0.17 | 0.82 | 0.69 | 0.50 | 0.26 | -0.42 | -0.15 | -0.11 | 0.07 | -0.37 |
| 32.5 | 302.5 | 0.63 | 0.42 | 0.33 | 0.35 | 0.46 | 0.21 | 0.65 | 0.48 | 0.34 | 0.27 | -0.64 | -0.27 | -0.40 | -0.13 | -0.59 |
| 32.5 | 307.5 | 0.69 | 0.43 | 0.48 | 0.54 | 0.69 | 0.20 | 0.46 | 0.33 | 0.25 | 0.26 | -0.63 | -0.15 | -0.37 | -0.12 | -0.54 |
| 32.5 | 312.5 | 0.80 | 0.51 | 0.44 | 0.44 | 0.68 | 0.26 | 0.45 | 0.41 | 0.26 | 0.28 | -0.28 | 0.08 | -0.16 | 0.05 | -0.21 |
| 32.5 | 317.5 | 0.83 | 0.47 | 0.16 | 0.26 | 0.61 | 0.36 | 0.47 | 0.49 | 0.14 | 0.10 | -0.01 | 0.21 | 0.04 | 0.23 | 0.24 |
| 32.5 | 322.5 | 0.62 | 0.16 | -0.19 | 0.10 | 0.44 | 0.39 | 0.41 | 0.43 | -0.04 | -0.09 | -0.10 | 0.10 | 0.09 | 0.33 | 0.45 |
| 32.5 | 327.5 | 0.24 | -0.29 | -0.54 | 0.05 | 0.27 | 0.29 | 0.21 | 0.23 | -0.35 | -0.19 | -0.21 | -0.03 | 0.09 | 0.40 | 0.52 |

Fig. 2.1 A subset of the monthly surface air temperature anomalies from the NOAAGlobalTemp Version 4.0 dataset.

1680 Following this spatial index as the row number, the climate data for a given

month is a column vector. If the dataset has data for p months, then the space-time data matrix has p columns. If the dataset has n grid boxes, then the data forms an $n \times p$ space-time data matrix. You can conveniently use row or column operations of a computer language to calculate statistics of the dataset, such as spatial average, temporal mean, temporal variance, etc.

For more explicit indication of space and time, you may use s for the row index and t for the column index in a space-time data matrix. Thus, $[A_{st}]$ indicates a space-time data matrix.

This space-time indexing can be extended to the data in 3D space and 1D time, as long as we can assign a unique ID s from 1 to n for a 3D grid box and a unique ID t for time. The presently popular netCDF (Network Common Data Form) data format in climate science, denoted by `.nc`, uses this index procedure for a 4D dataset. For example, to express the output of a 3D climate model, you can start your index longitude first for a given altitude and altitude. When longitude exhausts, count the next latitude. When the latitude exhausts, count the next altitude until the last layer of the atmosphere or ocean. Eventually, a space-time data matrix is formed $[A_{st}]$.

To visualize the row data of a space-time data matrix $[A_{st}]$, just plot a line graph of the row data against time. To visualize the column data of a space-time data matrix $[A_{st}]$, you need to convert the column vector into a 2D pixel format for a 2D domain (e.g., the contiguous United States (CONUS) region), or a 3D data array format for a 3D domain (e.g., the CONUS atmosphere domain from the 1,000 mb surface level to the 10 mb height level). This means that the climate data are represented in another matrix format, such as the 5° surface air temperature anomaly data for the entire world for December 2015 visualized by Fig. B.8. The data behind the figure is a 36×72 data matrix whose rows are for latitude and columns for longitude. This data matrix is in space-space pixel format like the data for a photo, and each time corresponds to a new space-space pixel data matrix. Thus, the latitude-longitude-time forms a data 3D array. With elevation, then latitude-longitude-altitude-time forms a 4D data array, which is often written in the netCDF file in climate science. You can use the 4DVD data visualization tool www.4dvd.org, described in Chapter 1, to visualize the 4D Reanalysis data array as an example to understand the space-time data plot, and netCDF data structure.

The coordinates (32.5, 262.5) in the 6th row of Fig. 2.1 indicates a $5^\circ \times 5^\circ$ grid box centered at (32.5°N, 97.5°W). This box covers part of Texas. The large temperature anomalies for the summer of 1934 (2.64°C for June, 2.12°C for July, and 2.36°C for August) were in the 1930s “Dust Bowl” period. The hot summer of 1934 was a wave of severe drought. The disastrous dust storms in the 1930s over the American and Canadian prairies destroyed many farms and greatly damaged the ecology.

2.2 Fundamental Properties of Matrices

This section provides a concise list to summarize fundamental properties and the commonly used operations of matrices. We limit our materials to the basics that are sufficient for this book.

- (i) Zero matrix: A zero matrix has its every entry equal to zero: $\mathbf{0} = [0]$, or $0 = [0]$, or explicitly

$$\mathbf{0} = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} \quad (2.2)$$

- (ii) Identity matrix: An identity matrix is a square matrix whose diagonal entries are all equal to one and whose off-diagonal entries are all equal to zero, and is denoted by I or \mathbf{I} . See an expression of an identity matrix below:

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad (2.3)$$

People also use the following notation

$$\mathbf{I} = [\delta_{ij}], \quad (2.4)$$

where

$$\delta_{ij} = \begin{cases} 1 & \text{when } i = j \\ 0 & \text{otherwise} \end{cases} \quad (2.5)$$

is called the Kronecker delta.

An identity matrix may be regarded as a special case of a *diagonal matrix*, which refers to any square matrix whose off-diagonal elements are zero. Hence, a diagonal matrix has the following general expression:

$$\mathbf{D} = [d_i \delta_{ij}], \quad (2.6)$$

where d_1, d_2, \dots, d_n are the n-diagonal elements. If $d_i = 1, i = 1, 2, \dots, n$, then the diagonal matrix becomes an identity matrix.

- (iii) A transpose matrix: The *transpose* of a matrix A is obtained by interchanging the rows and columns. The new matrix is denoted by \mathbf{A}^t . Computing the matrix transpose is very easy, simply rotating each horizontal row clockwise to a vertical column, one row at a time. If \mathbf{A} has dimension $n \times p$, its

1742 transpose has dimension $p \times n$. The elements of the transposed matrix are
1743 related to the originals by

$$(A^t)_{ij} = A_{ji} \quad (2.7)$$

1744 For example, if

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad (2.8)$$

1745 then

$$\mathbf{A}^t = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \quad (2.9)$$

1746 The equation $\mathbf{A}^t = \mathbf{A}$ or $a_{ij} = a_{ji}$ is equivalent to that the matrix \mathbf{A} is
1747 symmetric. Of course, a symmetric matrix must be a square matrix.

1748 (iv) Equal matrices: Two matrices \mathbf{A} and \mathbf{B} are equal if every pair of corresponding
1749 entries are equal, i.e., the equation $\mathbf{A} = \mathbf{B}$ is equivalent to $a_{ij} = b_{ij}$ for all
1750 i and j .

1751 (v) Matrix addition: The sum of two matrices \mathbf{A} and \mathbf{B} is defined by the sum of
1752 corresponding entries, i.e., $\mathbf{A} + \mathbf{B} = [a_{ij} + b_{ij}]$.

1753 (vi) Matrix subtraction: The difference of two matrices \mathbf{A} and \mathbf{B} is defined by the
1754 difference of corresponding entries, i.e., $\mathbf{A} - \mathbf{B} = [a_{ij} - b_{ij}]$. The equation
1755 of $\mathbf{A} = \mathbf{B}$ is equivalent to $\mathbf{A} - \mathbf{B} = 0$.

1756 (vii) Row vector: A row vector is of dimension p is a $1 \times p$ matrix:

$$\mathbf{u} = [u_1 \ u_2 \ \cdots \ u_p]. \quad (2.10)$$

1757 Matrix $\mathbf{A}_{n \times p}$ may be regarded as a stack of n row vectors $\mathbf{a}_{i:}, i = 1, 2, \dots, n$,
1758 each of which is a p -dimensional row vector. Hence,

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \vdots \\ \mathbf{a}_{n:} \end{bmatrix}. \quad (2.11)$$

1759 Here, $:$ in the second position of the double-index subscript means the
1760 inclusion of all the columns.

1761 (viii) Column vector: A column vector of dimension n is an $n \times 1$ matrix

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}. \quad (2.12)$$

1762 The transpose of a column vector becomes a row vector, and vice versa.

1763 Matrix $\mathbf{A}_{n \times p}$ may be regarded as an array of p column vectors $\mathbf{a}_{:,j}, j =$
1764 $1, 2, \dots, p$, each of which is an n -dimensional column vector. Hence,

$$\mathbf{A} = [\mathbf{a}_{:,1} \ \mathbf{a}_{:,2} \ \cdots \ \mathbf{a}_{:,p}]. \quad (2.13)$$

1765 Here, : in the first position of the double-index subscript means the inclusion
 1766 of all the rows.

1767 (ix) Dot product of two vectors: Two vectors of the same dimension can form a
 1768 *dot product* which is equal to the sum of the products of the corresponding
 1769 entries:

$$\mathbf{u} \cdot \mathbf{v} = u_1 v_1 + u_2 v_2 + \cdots + u_n v_n. \quad (2.14)$$

1770 The dot product is also called an *inner product*.

1771 As example, if

$$\mathbf{u} = [1 \ 2 \ 3], \quad \mathbf{v} = [4 \ 5 \ 6], \quad (2.15)$$

1772 then

$$\mathbf{u} \cdot \mathbf{v} = 1 \times 4 + 2 \times 5 + 3 \times 6 = 32. \quad (2.16)$$

1773 The *amplitude* of vector \mathbf{u} of dimension n is defined as

$$|\mathbf{u}| = \sqrt{u_1^2 + u_2^2 + \cdots + u_n^2}. \quad (2.17)$$

1774 Sometimes, the amplitude is also called length, or Euclidean length, or
 1775 magnitude. Please do not mix the concept of Euclidean length of a vector
 1776 with the dimensional length of a vector. The latter means the number of
 1777 entries of a vector, i.e., n .

1778 If the Euclidean length of \mathbf{u} is equal to one, we say that the \mathbf{u} is a *unit*
 1779 *vector*. If every element of \mathbf{u} is zero, then we say that \mathbf{u} is a *zero vector*.

1780 By the definition of dot product, we have

$$|\mathbf{u}|^2 = \mathbf{u} \cdot \mathbf{u}. \quad (2.18)$$

1781 If $\mathbf{u} \cdot \mathbf{v} = 0$, we say that \mathbf{u} and \mathbf{v} are *orthogonal*. Further, if $\mathbf{u} \cdot \mathbf{v} = 0$ and
 1782 $|\mathbf{u}| = |\mathbf{v}| = 1$, then we say that \mathbf{u} and \mathbf{v} are *orthonormal*.

1783 (x) Matrix multiplication: The product of matrix $\mathbf{A}_{n \times p}$ and matrix $\mathbf{B}_{p \times m}$ is an
 1784 $n \times m$ matrix $\mathbf{C}_{n \times m}$ whose element c_{ij} is the dot product of the i th row
 1785 vector of A and j th column vector of B :

$$c_{ij} = \mathbf{a}_{i:} \cdot \mathbf{b}_{:j}. \quad (2.19)$$

1786 We denote

$$\mathbf{C}_{n \times m} = \mathbf{A}_{n \times p} \mathbf{B}_{p \times m}, \quad (2.20)$$

1787 or

$$\mathbf{C} = \mathbf{AB}. \quad (2.21)$$

1788 Note that the number of columns of \mathbf{A} and that of rows of \mathbf{B} must be the
 1789 same when the multiplication \mathbf{AB} can be made, because the dot product
 1790 $\mathbf{a}_{i:} \cdot \mathbf{b}_{:j}$ requires this condition. This is referred to as the dimension-matching

1791 condition for matrix multiplication. If this condition is violated, the two
 1792 matrices cannot be multiplied. For example, for the following two matrices

$$\mathbf{A}_{3 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 4 \\ 3 & 2 \end{bmatrix} \quad \mathbf{B}_{2 \times 2} = \begin{bmatrix} 0 & -1 \\ 1 & 2 \end{bmatrix} \quad (2.22)$$

1793 We can compute

$$\mathbf{A}_{3 \times 2} \mathbf{B}_{2 \times 2} = \begin{bmatrix} 0 & -1 \\ 4 & 8 \\ 2 & 1 \end{bmatrix} \quad (2.23)$$

1794 However, the expression

$$\mathbf{B}_{2 \times 2} \mathbf{A}_{3 \times 2} \quad (2.24)$$

1795 is not defined, because the dimensions do not match. Thus, for matrix mul-
 1796 tiplication of two matrices, their order is important. The product \mathbf{BA} may
 1797 not be equal to \mathbf{AB} even when both are defined. That is, the commutative
 1798 law does not hold for matrix multiplication.

1799 The dot product of two vectors can be written as the product of two
 1800 matrices. If \mathbf{u} and \mathbf{v} are n-dimensional column vectors, then

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^t \mathbf{v}. \quad (2.25)$$

1801 The right hand side is a $1 \times n$ matrix times an $n \times 1$ matrix, and the product
 1802 is a 1×1 matrix, whose element is the result of the dot product. Computer
 1803 usually programs dot product in this way of matrix multiplication.

1804 A scalar can always multiply a matrix, which is defined as follows. Given
 1805 a scalar c and a matrix \mathbf{A} , their product is

$$c\mathbf{A} = [ca_{ij}] = \mathbf{Ac}. \quad (2.26)$$

1806 The scalar multiplication can be extended to multiple vectors
 1807 $(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p)$
 1808 or matrices to form a *linear combination*:

$$\mathbf{u} = c_1 \mathbf{u}_1 + c_2 \mathbf{u}_2 + \dots + c_p \mathbf{u}_p. \quad (2.27)$$

1809 where c_1, c_2, \dots, c_p are coefficients of the linear combination and at least
 1810 one of the coefficients is non-zero. Multivariate linear regression discussed
 1811 at the end of last chapter is a linear combination. This is a very useful
 1812 mathematical expression in data science.

1813 (xi) Matrix inversion: For a given square matrix \mathbf{A} , if there is a matrix \mathbf{B} such
 1814 that

$$\mathbf{BA} = \mathbf{AB} = \mathbf{I}, \quad (2.28)$$

1815 then \mathbf{B} is called the *inverse matrix* of \mathbf{A} , denoted by \mathbf{A}^{-1} , i.e.,

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{AA}^{-1} = \mathbf{I}. \quad (2.29)$$

1816 Not all the matrices have an inverse. If a matrix has an inverse, then the
1817 matrix is said to be *invertible*. Equivalently, \mathbf{A}^{-1} exists.

1818 As an example of the matrix inversion, given

$$\mathbf{A} = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}, \quad (2.30)$$

1819 we have

$$\mathbf{A}^{-1} = \begin{bmatrix} 2/3 & 1/3 \\ -1/3 & 1/3 \end{bmatrix}. \quad (2.31)$$

1820 Hand-calculation for the inverse of a small matrix is already very difficult,
1821 and that for the inverse of a large matrix is almost impossible. Computers
1822 can do the calculations for us, as shown in examples later.

1823 According to the definition of inverse, we have the following formula for
1824 the inverse of the product of two matrices:

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}, \quad (2.32)$$

1825 if both \mathbf{A} and \mathbf{B} are invertible matrices. Please note the order switch of
1826 the matrices.

1827 With the definition of an inverse matrix, we can define the *matrix division*
1828 by

$$\mathbf{A}/\mathbf{B} = \mathbf{AB}^{-1} \quad (2.33)$$

1829 when \mathbf{B}^{-1} exists. In matrix operations, we usually do not use the concept
1830 of matrix division, but always use the matrix inverse and matrix multipli-
1831 cation.

1832 (xii) More properties of the matrix transpose:

$$(\mathbf{A}^t)^t = \mathbf{A} \quad (2.34)$$

$$(\mathbf{A} + \mathbf{B})^t = \mathbf{A}^t + \mathbf{B}^t \quad (2.35)$$

$$(\mathbf{AB})^t = \mathbf{B}^t \mathbf{A}^t \quad (2.36)$$

$$(\mathbf{A}^{-1})^t = (\mathbf{A}^t)^{-1}. \quad (2.37)$$

1833 (xiii) Orthogonal matrices: An *orthogonal matrix*¹ is one whose row vectors are
1834 orthonormal. In this case, the inverse matrix can be easily found: It is its
1835 transpose. That is, if \mathbf{A} is an orthogonal matrix, then

$$\mathbf{A}^{-1} = \mathbf{A}^t. \quad (2.38)$$

1836 The proof of this claim is very simple. The orthonormal property of the
1837 row vectors of \mathbf{A} implies that

$$\mathbf{AA}^t = \mathbf{I}. \quad (2.39)$$

1838 By the definition of matrix inverse, \mathbf{A}^t is the inverse matrix of \mathbf{A} .

¹ Although *orthogonal matrix* is a standard mathematical terminology, it is acceptable if you call it *orthonormal matrix*

1839 If \mathbf{A} is an orthogonal matrix, its row vectors are also orthonormal. This
 1840 can be proved by multiplying both sides of the above by \mathbf{A}^t from the left:

$$\mathbf{A}^t \mathbf{A} \mathbf{A}^t = \mathbf{A}^t \mathbf{I}. \quad (2.40)$$

1841 Then multiply both sides of this equation by $(\mathbf{A}^t)^{-1}$ from the right:

$$\mathbf{A}^t \mathbf{A} (\mathbf{A}^t (\mathbf{A}^t)^{-1}) = \mathbf{A}^t \mathbf{I} (\mathbf{A}^t)^{-1}, \quad (2.41)$$

1842 which yields

$$\mathbf{A}^t \mathbf{A} = \mathbf{I}. \quad (2.42)$$

1843 This implies that the column vectors of \mathbf{A} are orthonormal.

1844 As an example, the following matrix

$$\mathbf{T} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad (2.43)$$

1845 is an orthogonal matrix for any given real number θ . You can easily verify
 1846 this using the trigonometrical identity $\sin^2 \theta + \cos^2 \theta = 1$. We thus have

$$\mathbf{T}^{-1} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}, \quad (2.44)$$

1847 You can easily verify that $\mathbf{T}^{-1} \mathbf{T} = \mathbf{I}$ by hand-calculation of the product
 1848 of the two matrices in this equation.

1849 2.3 Some basic concepts and theories of linear algebra

1850 According to Encyclopedia.com, “linear algebra originated as the study of linear
 1851 equations.” Linear algebra deals with vectors, matrices and vector spaces. Before
 1852 the 1950s, it was part of Abstract Algebra (Tucker 1993). In 1965 the Commit-
 1853 tee on the Undergraduate Program in Mathematics, Mathematical Association of
 1854 America, outlined the following topics for a stand-alone linear algebra course: Lin-
 1855 ear systems, matrices, vectors, linear transformations, unitary geometry with char-
 1856 acteristic values. The vectors and matrices have been dealt in the previous two
 1857 sections of this chapter. This section deals with linear systems of equations, and
 1858 linear transformations, and next with characteristic values.

1860 2.3.1 Linear equations

1861 A meteorologist needs to make a decision on what instruments to order under the
 1862 following constraint. She is given a budget of \$10,000 to purchase 30 instruments for
 1863 her observational sites. Her supplier has two products for the instrument: The first
 1864 is \$30 per set, and the second \$40 per set. She would like to buy the second type of
 1865 instrument as many as possible under the budget constraint. Then, the question is

1867 how many instruments of the second kind she can buy? This problem leads to the
1868 following linear system of two equations:

$$30x_1 + 40x_2 = 1000, \quad (2.45)$$

$$x_1 + x_2 = 30. \quad (2.46)$$

1869 The solution to these linear equations is $x_1 = 20$ and $x_2 = 10$.

1870 This system of linear equations can be expressed in matrix and vectors as follows:

$$\mathbf{Ax} = \mathbf{b}, \quad (2.47)$$

1871 where

$$\mathbf{A} = \begin{bmatrix} 30 & 40 \\ 1 & 1 \end{bmatrix} \quad (2.48)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (2.49)$$

$$\mathbf{b} = \begin{bmatrix} 1000 \\ 30 \end{bmatrix} \quad (2.50)$$

$$(2.51)$$

1872 Then, the solution of this system may be tightly expressed in the following way:

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}. \quad (2.52)$$

1873 This expression is convenient for mathematical proofs, but is rarely used for solving
1874 a linear system, because finding an inverse matrix is computationally costly. A way
1875 to solve a linear system is to use Gauss elimination. The corresponding computing
1876 procedure is called the row operation on a matrix. There are numerous ways of
1877 solving a linear system. Some are particularly efficient for a certain system, such
1878 as a sparse matrix or a matrix of a diagonal bend of width equal 3 or 5. Efficient
1879 algorithms for a linear system, particularly an extremely large system, are forever
1880 a research topic. In this book we use a computer to solve a linear system without
1881 studying the algorithm details. The R and Python commands are below:

```
1882 solve(A, b) #R code for finding x
1883 numpy.linalg.solve(A, b) #Python code
```

1884 2.3.2 Linear transformations

1886 A *linear transformation* is to convert vector $\mathbf{x}_{n \times 1}$ into $\mathbf{y}_{m \times 1}$ using the multiplication
1887 of a matrix $\mathbf{T}_{m \times n}$:

$$\mathbf{y}_{m \times 1} = \mathbf{T}_{m \times n} \mathbf{x}_{n \times 1} \quad (2.53)$$

1888 For example, the matrix

$$\mathbf{T} = \begin{bmatrix} -0.1 & 4 \\ 0.1 & -3 \end{bmatrix} \quad (2.54)$$

1889 transforms the vector

$$\begin{bmatrix} 1000 \\ 30 \end{bmatrix} \quad (2.55)$$

1890 into

$$\begin{bmatrix} 20 \\ 10 \end{bmatrix} \quad (2.56)$$

1891 This is the solution of the linear system in the previous subsection.

1892 Usually, the linear transformation $\mathbf{T}\mathbf{x}$ changes both direction and magnitude of
1893 \mathbf{x} . However, if \mathbf{T} is an orthogonal matrix, then $\mathbf{T}\mathbf{x}$ does not change the magnitude of
1894 \mathbf{x} , and changes only the direction. This claim can be simply proved by the following
1895 formula:

$$|\mathbf{T}\mathbf{x}|^2 = (\mathbf{T}\mathbf{x})^t \mathbf{T}\mathbf{x} = \mathbf{x}^t \mathbf{T}^t \mathbf{T}\mathbf{x} = \mathbf{x}^t (\mathbf{T}^t \mathbf{T})\mathbf{x} = \mathbf{x}^t \mathbf{I}\mathbf{x} = |\mathbf{x}|^2. \quad (2.57)$$

1896 Thus, if \mathbf{T} is an orthogonal matrix, then $\mathbf{T}\mathbf{x}$ is a rotation of the vector \mathbf{x} .

1897 2.3.3 Linear independence

1899 The vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ are *linearly dependent* if no vector can be represented
1900 by a linear combination of other $p - 1$ vectors in this group. Otherwise, the group
1901 of vectors are linearly independent.

1902 If it is not linearly independent, then there must be a vector which can be rep-
1903 resented by the other vectors through a linear combination. Suppose this vector is
1904 \mathbf{x}_1 , then

$$\mathbf{x}_1 = d_2 \mathbf{x}_2 + \dots + d_p \mathbf{x}_p, \quad (2.58)$$

1905 where at least one of the coefficients d_2, d_3, \dots, d_p is non-zero. Thus, the linear
1906 system of equations for $c_1, c_2, c_3, \dots, c_p$

$$c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \dots + c_p \mathbf{x}_p = 0 \quad (2.59)$$

1907 has a non-zero solution. This system can be written as a matrix form

$$\mathbf{X}\mathbf{c} = \mathbf{0}, \quad (2.60)$$

1908 where column vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ form the matrix \mathbf{X}

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_p], \quad (2.61)$$

1909 the unknown vector is \mathbf{c}

$$\mathbf{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix}, \quad (2.62)$$

1910 and $\mathbf{0}$ is the p -dimensional zero column vector. The solution of this matrix equation
1911 is

$$\mathbf{c} = \mathbf{X}^{-1} \mathbf{0} = \mathbf{0}. \quad (2.63)$$

1912 However, \mathbf{c} must not be zero. This contradiction implies that \mathbf{X}^{-1} does not exist if
1913 the column vectors are linearly dependent. In other words, if \mathbf{X}^{-1} exists, then
1914 its column vectors are linearly independent.

1915 Consider vectors in a three dimensional space. Any two column vectors \mathbf{x}_2 and
1916 \mathbf{x}_3 define a plane. If \mathbf{x}_1 can be written as a linear combination of \mathbf{x}_2 and \mathbf{x}_3 , then
1917 it must lie in the same plane. So, \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{x}_3 are linearly dependent. The matrix
1918 $[\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_3]_{3 \times 3}$ is not invertible.

1919 2.3.4 Determinants

1921 For a square matrix \mathbf{A} , a convenient notation and concept is its *determinant*. It is
1922 a scalar and is denoted by $\det(\mathbf{A})$ or $|\mathbf{A}|$. For a 2×2 matrix

$$\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, \quad (2.64)$$

1923 its determinant is

$$\det[\mathbf{A}] = ad - cd. \quad (2.65)$$

1924 For a higher-dimensional matrix, the computation is quite complex and is computationally expensive. We usually do not need to calculate the determinant of a large
1925 matrix, say $\mathbf{A}_{172 \times 172}$. The computer command for computing the determinant of a
1926 small square matrix is as follows:

```
1928 det(A) #R command for determinant
1929 np.linalg.det(a) #Python command for determinant
```

1930 Two 2-dimensional column vectors \mathbf{x}_1 and \mathbf{x}_2 can span a parallelogram, whose
1931 area S is equal to the absolute value of the determinant of the matrix consisting of
1932 the two vectors $\mathbf{A} = [\mathbf{x}_1 \ \mathbf{x}_2]$:

$$S = |\det[\mathbf{x}_1 \ \mathbf{x}_2]|. \quad (2.66)$$

1933 Three 3-dimensional column vectors $\mathbf{x}_1, \mathbf{x}_2$ and \mathbf{x}_3 can span a parallelepiped,
1934 whose volume V is equal to the absolute value of the determinant of the matrix
1935 consisting of the three vectors $\mathbf{A} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_3]$:

$$V = |\det[\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_3]|. \quad (2.67)$$

1936 A few commonly used properties of determinant are listed below

- 1937 (a) The determinant of a diagonal matrix is the product of its diagonal elements.
- 1938 (b) If a determinant has a zero row or column, the determinant is zero.
- 1939 (c) The determinant does not change after a matrix transpose, i.e., $\det[\mathbf{A}^t] = \det[\mathbf{A}]$.
- 1941 (d) The determinant of the product of two matrices: $\det[\mathbf{AB}] = \det[\mathbf{A}]\det[\mathbf{B}]$.
- 1942 (e) The determinant of the product of a matrix with a scalar: $\det[c\mathbf{B}] = c^n\det[\mathbf{B}]$,
1943 if \mathbf{A} is an $n \times n$ matrix.
- 1944 (f) The determinant of an orthogonal matrix is equal to 1 or -1.

1945
1946

2.3.5 Rank of a matrix

1947 The *rank* of \mathbf{A} is the greatest number of columns of the matrix that are linearly
 1948 independent, and is denoted by $r[\mathbf{A}]$.

1949 For a 3×3 matrix \mathbf{A} , if we treat each column as a 3-dimensional vector. If all
 1950 three lie along a line (i.e., collinear), the rank of \mathbf{A} is one. If all three of the vectors
 1951 lie in a plane, but are not collinear, the rank is two. If none of the three are collinear
 1952 or lie in a plane, the rank is three.

1953 If the rank of a square matrix is less than its dimension, then at least one column
 1954 can be a linear combinations of other columns, which implies that the determinant
 1955 vanishes. If \mathbf{A} has rank r , it is possible to find r linearly independent columns, and
 1956 all the other columns are linear combinations of these r independent columns.

1957 Some properties about the matrix rank are listed below:

- 1958 (a) If $\det[\mathbf{A}_{n \times n}] \neq 0$, then $r[\mathbf{A}] = n$, and the matrix $\mathbf{A}_{n \times n}$ is invertible and is
 1959 said to be *nonsingular*.
- 1960 (b) If $\det[\mathbf{A}] = 0$, then the rank of \mathbf{A} is less than n , and \mathbf{A} is not invertible and
 1961 is said to be *singular*.
- 1962 (c) If \mathbf{B} is multiplied by a nonsingular matrix \mathbf{A} , the product has the same rank
 1963 as \mathbf{B} .
- 1964 (d) $0 \leq r[\mathbf{A}_{n \times p}] \leq \min(n, p)$.
- 1965 (e) $r[\mathbf{A}] = r[\mathbf{A}^t]$.
- 1966 (f) $r[\mathbf{AB}] \leq \min(r[\mathbf{A}], r[\mathbf{B}])$.
- 1967 (g) $r[\mathbf{AA}^t] = r[\mathbf{A}^t \mathbf{A}] = r[\mathbf{A}]$.
- 1968 (h) $r[\mathbf{A} + \mathbf{B}] \leq r[\mathbf{A}] + r[\mathbf{B}]$.

1969 Computers can easily demonstrate the matrix computations following the theories
 1970 presented in this chapter so far. The computer code is below.

```

1971 #R code: Computational examples of matrices
1972 A = matrix(c(1,0,0,4,3, 2), nrow = 3, byrow = TRUE)
1973 B = matrix(c(0,1,-1,2), nrow = 2) #form a matrix by columns
1974 C = A%*%B #matrix multiplication
1975 C
1976 #[1,] 0 -1
1977 #[2,] 4 8
1978 #[3,] 2 1
1979 t(C) # transpose matrix of C
1980 #[1,] 0 4 2
1981 #[2,] -1 8 1
1982
1983 A = matrix(c(1, -1, 1, 2), nrow =2, byrow = TRUE)
1984 solve(A) #compute the inverse of A
1985 #[1,] 0.6666667 0.3333333
1986 #[2,] -0.3333333 0.3333333
1987 A%*%solve(A) #verify the inverse of A
1988 #[1,] 1.000000e+00 0
1989 #[2,] 1.110223e-16 1
1990
1991 #Solve linear equations

```

```
1992 A = matrix(c(30, 40, 1, 1), nrow = 2, byrow = TRUE)
1993 b = c(1000, 30)
1994 solve(A,b)
1995 #[1] 20 10
1996 solve(A) %*% b #Another way to solve the equations
1997 det(A) #compute the determinant
1998 #[1] -10
1999
2000 library(Matrix)
2001 rankMatrix(A) #Find the rank of a matrix
2002 #[1] 2 #rank(A) = 2
2003
2004 #Orthogonal matrices
2005 p = sqrt(2)/2
2006 Q = matrix(c(p,-p,p,p), nrow=2)
2007 Q #is an orthogonal matrix
2008 # [,1] [,2]
2009 #[1,] 0.7071068 0.7071068
2010 #[2,] -0.7071068 0.7071068
2011 Q %*% t(Q) #verify 0 as an orthogonal matrix
2012 # [,1] [,2]
2013 #[1,] 1 0
2014 #[2,] 0 1
2015 det(Q) #The determinant of an orthogonal matrix is 1 or -1
2016 #[1] 1
```

```

# Python matrix multiplication
A = [[1, 0], [0, 4], [3, 2]]
B = [[0,-1],[1,2]]
C = np.matmul(A,B) #Or C = np.dot(A,B)
print('C=', C)
#C= [[ 0 -1]
# [ 4  8]
# [ 2  1]]
print('Transpose\u20d7matrix\u20d7of\u20d7C\u20d7= ', C.transpose())
#Transpose matrix of C = [[ 0  4  2]
# [-1  8  1]]

#matrix inversion
A = [[1,-1],[1,2]]
np.linalg.inv(A)# compute the inverse of A
#array([[ 0.66666667,  0.33333333],
#       [-0.33333333,  0.33333333]])

#Solve a system of linear equations
A = [[30, 40],[1, 1]]
b = [[1000],[30]]
x = np.linalg.solve(A,b)
print('x=', x)
#x= [[20.]
# [10.]]

#Compute determinant of the previous matrix A
print('Determinant\u20d7det(A)= ', np.linalg.det(A))
#Determinant det(A)= -10.000000000000002

#An orthogonal matrix
p = np.sqrt(2)/2
Q = [[p,p],[-p,p]]
print('Orthogonal\u20d7matrix\u20d7Q=\u20d7, np.round(Q,2)')
T = np.transpose(Q)
print('Q\u20d7times\u20d7transpose\u20d7of\u20d7Q\u20d7=\u20d7, np.matmul(Q,T))')
print('Determinant\u20d7of\u20d7Q\u20d7=\u20d7, np.linalg.det(Q))')
#Orthogonal matrix Q= [[ 0.71  0.71]
# [-0.71  0.71]]
#Q times transpose of Q =  [[1.  0.]
# [0.  1.]]
#Determinant of Q = 1.0

```

2018

2.4 Eigenvectors and eigenvalues

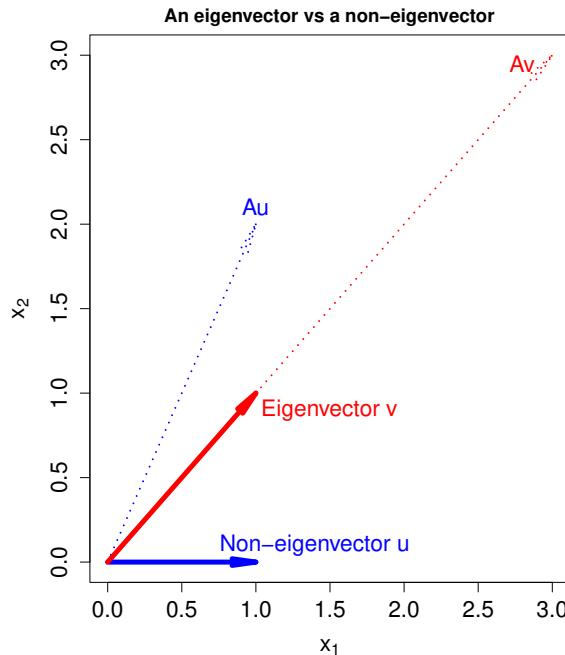
2019

2020
2021

2.4.1 Definition of eigenvectors and eigenvalues

2022 The linear transform $\mathbf{y} = \mathbf{A}\mathbf{u}$ usually results in \mathbf{y} not parallel to \mathbf{u} . For example,

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad (2.68)$$

2023 The vectors $\mathbf{u} = (1, 0)$ and $\mathbf{A}\mathbf{u} = (1, 2)$ are not parallel in the 2-dimensional space.
2024 See the blue vectors in Fig. 2.2.**Fig. 2.2**

An eigenvector \mathbf{v} , a non-eigenvector \mathbf{u} , and their linear transforms by matrix \mathbf{A} : \mathbf{Av} and \mathbf{Au} . Here, \mathbf{Av} and \mathbf{v} are parallel, and \mathbf{Au} and \mathbf{u} are not parallel.

2025 However, there exist some special vectors \mathbf{v} such that \mathbf{Av} is parallel to \mathbf{v} .
2026 For example, $\mathbf{v} = (1, 1)$ is such a vector, since $\mathbf{Av} = (3, 3)$ is in the same direction
2027 as $\mathbf{v} = (1, 1)$. See the red vectors in Fig. 2.2. If two vectors are parallel, then one
2028 vector is a scalar multiplication of the other, e.g., $(3, 3) = 3(1, 1)$. We denote this
2029 scalar by λ . Thus,

$$\mathbf{Av} = \lambda\mathbf{v}. \quad (2.69)$$

2030 These vectors v are special to A , maintain their own orientation when multiplied
2031 by A , and are called *eigenvectors*. Here, “eigen” is from German, meaning “self”,

2032 “own”, “particular”, or “special.”² The corresponding scalars λ are called *eigenvalues*.
 2033 The formula (2.69) is a mathematical definition of the eigenvalue problem
 2034 for matrix \mathbf{A} .

2035 If \mathbf{v} is an eigenvector, then its multiplication to a scalar c is also an eigenvector,
 2036 since

$$\mathbf{A}(c\mathbf{v}) = c\mathbf{A}\mathbf{v} = c\lambda\mathbf{v} = \lambda(c\mathbf{v}).$$

2037 Namely, all the vectors in the same direction \mathbf{v} are also eigenvectors. The eigenvectors
 2038 of length equal one are called unit eigenvectors, or unitary eigenvectors, and
 2039 are unique up to a positive or negative sign. Most computer programs output unit
 2040 eigenvectors. Thus, eigenvector describes a direction or an orientation. Each square
 2041 matrix has its own special orientations.

2042 The aforementioned vector

$$\mathbf{v} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad (2.70)$$

2043 is en eigenvector that maintains its own direction after multiplied by \mathbf{A} :

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 1 \end{bmatrix}. \quad (2.71)$$

2044 Thus,

$$\mathbf{v} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

2045 is an eigenvector of A and $\lambda = 3$ is an eigenvalue of A . The corresponding unit
 2046 eigenvector is

$$\mathbf{e} = \mathbf{v}/|\mathbf{v}| = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

2047 Another eigenvector for the above matrix \mathbf{A} is $\mathbf{v}_2 = (1, -1)$:

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ -1 \end{bmatrix} = -1 \times \begin{bmatrix} 1 \\ -1 \end{bmatrix}. \quad (2.72)$$

2048 The second eigenvalue is $\lambda_2 = -1$.

2049 The computer code for the eigenvalues and eigenvectors of the above matrix \mathbf{A}
 2050 is below.

```
2051 #R code for eigenvectors and eigenvalues
2052 A = matrix(c(1, 2, 2, 1), nrow=2)
2053 eigen(A)
```

² The “eigen” part in the word “eigenvector” is from German or Dutch and means “self” or “own”, as in “one’s own.” Thus, an eigenvector v is A ’s “own” vector. In English books, the word “eigenvector” is the standard translation of the German word “eigenvektor.” The word “eigenvalue” is translated from the German word “eigenwert”, as “wert” means “value.” Instead of eigenvector, some English publications use “characteristic vector”, which indicates “characteristics of a matrix”, or “its own property of a matrix.” German mathematician David Hilbert (1862-1943) was the first to use “eigenvektor”, and “eigenwert” in his 1904 article about a general theory of linear integral equations.

```

2054 #$values
2055 #[1] 3 -1
2056 #$vectors
2057 #[,1]      [,2]
2058 #[1,] 0.7071068 -0.7071068
2059 #[2,] 0.7071068 0.7071068

```

```

#Python code for eigenvectors and eigenvalues
A = [[1,2], [2,1]]
np.linalg.eig(A)
#(array([ 3., -1.]), array([[ 0.70710678, -0.70710678],
#                           [ 0.70710678,  0.70710678]]))

```

2060
2061 If \mathbf{A} is an $N \times N$ matrix, then it has N eigenvalues and eigenvectors $(\lambda_n, \mathbf{v}_n)$, $n = 1, 2, \dots, N$ for the following reason. The eigenvector \mathbf{v} satisfies

$$(2.73) \quad (\mathbf{A} - \lambda \mathbf{I})\mathbf{v} = 0.$$

2063 The non-zero solution \mathbf{v} of this equation requires that

$$(2.74) \quad \det(\mathbf{A} - \lambda \mathbf{I}) = 0.$$

2064 Expanding the determinant out leads to an N th degree polynomial in λ , which has
2065 exactly N roots. Some roots may be repeated and hence counted multiple times
2066 toward N . Some roots may be complex numbers, which are not discussed in this
2067 book. Each root is an eigenvalue and corresponds to an eigenvector.

2068 The above concise determinant expression is tidy and useful for mathematical
2069 proofs, but is not used as a computer algorithm for calculating eigenvalues or eigen-
2070 vectors, because it is computationally costly or even impossible for a large matrix.
2071 Nonetheless, some traditional textbooks of linear algebra defined eigenvalues using
2072 Eq. (2.74). This traditional definition of eigenvalues appears to be abstract. The
2073 eigenvector equation (2.73) seems not providing an image of the eigenvector. Thus,
2074 many students forget the definition of eigenvalues and eigenvectors shortly after the
2075 final exams of a linear algebra course.

2076 Figure 2.2 may be generated by the following computer code.

```

2077 #R plot eigenvector v vs a non-eigenvector u
2078 #Create your working directory named LinAlg and go there
2079 setwd('/Users/sshenn/LinAlg')
2080 getwd() #Verify that you are in the directory/folder
2081 #[1] "/Users/sshenn/LinAlg"
2082
2083 setEPS() #Plot the figure and save the file
2084 postscript("fig0502.eps", width = 6)
2085 par(mar=c(4.5,4.5,2.0,0.5))
2086 plot(9,9,
2087       main = 'An_eigenvector_vs_a_non-eigenvector',
2088       cex.axis = 1.4, cex.lab = 1.4,
2089       xlim = c(0,3), ylim=c(0,3),
2090       xlab = bquote(x[1]), ylab = bquote(x[2]))

```

```

2091 arrows(0,0, 1,0, length = 0.25,
2092         angle = 8, lwd = 5, col = 'blue')
2093 arrows(0,0, 1,2, length = 0.3,
2094         angle = 8, lwd = 2, col = 'blue', lty = 3)
2095 arrows(0,0, 1,1, length = 0.25,
2096         angle = 8, lwd = 5, col='red')
2097 arrows(0,0, 3,3, length = 0.3,
2098         angle = 8, lwd = 2, col='red', lty = 3)
2099 text(1.4,0.1, 'Non-eigenvector $\omega$  $u$ ', cex =1.4, col = 'blue')
2100 text(1.0,2.1, 'Au', cex =1.4, col = 'blue')
2101 text(1.5,0.9, 'Eigenvector $\omega$  $v$ ', cex =1.4, col = 'red')
2102 text(2.8, 2.95, 'Av', cex =1.4, col = 'red')
2103 dev.off()

```

```

#Python plot eigenvector vs a non-eigenvector
import matplotlib.patches as patches
fig = plt.figure(figsize = (12,12))

plt.axes().set_xlim(-0.1,3.1)
plt.axes().set_ylim(-0.1,3.1)
plt.axes().set_aspect(1)
style = "Simple, tail_width=0.5, head_width=8, head_length=18"
kw1 = dict(arrowstyle=style, color="blue")
kw2 = dict(arrowstyle=style, color="red")

a1 = patches.FancyArrowPatch((0,0), (1,0), **kw1,
                             linewidth =5)
a2 = patches.FancyArrowPatch((0, 0), (1,2),**kw1)
a3 = patches.FancyArrowPatch((0, 0), (1,1), **kw2,
                             linewidth = 5)
a4 = patches.FancyArrowPatch((0, 0), (3,3),**kw2)
for a in [a1, a2, a3, a4]:
    plt.gca().add_patch(a)
plt.title('An eigenvector  $\omega v$  vs a non-eigenvector  $\omega u$ ')
plt.xlabel(r'$x_1$', fontsize = 25)
plt.ylabel(r'$x_2$', fontsize = 25)
plt.text(0.6, 0.1, 'Non-eigenvector $\omega$  $u$ ',
         color = 'blue', fontsize =25)
plt.text(0.8, 2.0, 'Au',
         color = 'blue', fontsize =25)
plt.text(1.03,0.85, 'Eigenvector $\omega$  $v$ ',
         color = 'red', fontsize =25)
plt.text(2.7, 2.9, 'Av',
         color = 'red', fontsize =25)
plt.show()

```

2105 **2.4.2 Properties of eigenvectors and eigenvalues for a symmetric
2106 matrix**
2107

2108 A covariance matrix or a correlation matrix is a symmetric matrix and is often used
2109 in climate science. For a symmetric matrix \mathbf{A} , its eigenvalues and eigenvectors have
2110 the following properties:

- 2111 (a) Eigenvalues of a symmetric matrix are real numbers.
- 2112 (b) The n different unit eigenvectors of a symmetric matrix $\mathbf{A}_{n \times n}$ are independent
2113 and form an orthonormal set, i.e., $\mathbf{e}^{(\ell')} \cdot \mathbf{e}^{(\ell)} = \delta_{\ell\ell'}$, where $\mathbf{e}^{(\ell')}$ and $\mathbf{e}^{(\ell)}$ are
2114 any two different unit eigenvectors. We can use these unit vectors to express
2115 any n -dimensional vector using a linear combination:

$$\mathbf{x} = c_1 \mathbf{e}^{(1)} + c_2 \mathbf{e}^{(2)} + \cdots + c_n \mathbf{e}^{(n)}, \quad (2.75)$$

2116 where c_1, c_2, \dots, c_n are scalar coefficients of the linear combination.

- 2117 (c) If all the eigenvalues of $\mathbf{A}_{n \times n}$ are positive, then the *quadratic form*

$$Q(\mathbf{x}) = \mathbf{x}^t \mathbf{A} \mathbf{x} = \sum_{i,j=1}^n a_{ij} x_i x_j \quad (2.76)$$

2118 is also a positive scalar for any non-zero vector \mathbf{x} . A quadratic form may
2119 be used to express kinetic energy or total variance in climate science. The
2120 kinetic energy in climate science is always positive. For all the unit vectors
2121 \mathbf{x} , the maximum of the quadratic form is equal to the largest eigenvalue of
2122 \mathbf{A} . The maximum is achieved when \mathbf{x} is the corresponding eigenvector of
2123 \mathbf{A} .

- 2124 (d) The rank of matrix \mathbf{A} is equal to the number of nonzero eigenvalues, where
2125 the multiplicity of repeated eigenvalues is counted.
- 2126 (e) Eigenvalues of a diagonal matrix are equal to the diagonal elements.
- 2127 (f) For a symmetric matrix $\mathbf{A}_{n \times n}$, its n unit column eigenvectors form an or-
2128 thogonal matrix $\mathbf{Q}_{n \times n}$ such that $\mathbf{A}_{n \times n}$ can be diagonalized by $\mathbf{Q}_{n \times n}$ in
2129 the following way

$$\mathbf{Q}_{n \times n}^t \mathbf{A}_{n \times n} \mathbf{Q}_{n \times n} = \mathbf{D}, \quad (2.77)$$

2130 where \mathbf{D} is a diagonal matrix whose diagonal elements are eigenvalues
2131 $\lambda_1, \lambda_2, \dots, \lambda_n$ of \mathbf{A} . Further, the column vectors of \mathbf{Q} are the unit eigen-
2132 vectors of \mathbf{A} . This is a matrix diagonalization process.

2133 For the symmetric matrix \mathbf{A} in Eq. (4.51)

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, \quad (2.78)$$

2134 its eigenvalues and eigenvectors are

$$\lambda_1 = 3, \quad \lambda_2 = -1, \quad (2.79)$$

$$\mathbf{v}_1 = \begin{bmatrix} \sqrt{2}/2 \\ \sqrt{2}/2 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} -\sqrt{2}/2 \\ \sqrt{2}/2 \end{bmatrix} \quad (2.80)$$

2135 Thus, the orthogonal matrix \mathbf{A} and the diagonal matrix \mathbf{d} are as follows

$$\mathbf{Q} = \begin{bmatrix} \sqrt{2}/2 & \sqrt{2}/2 \\ \sqrt{2}/2 & -\sqrt{2}/2 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix} \quad (2.81)$$

2136 With \mathbf{Q} , \mathbf{A} , and \mathbf{D} , you can easily verify Eq. (2.77) by hand calculation or by
2137 computer coding.

2138 Equation (2.77) can be written in the following way

$$\mathbf{A}_{n \times n} = \mathbf{Q}_{n \times n} \mathbf{D}_{n \times n} \mathbf{Q}_{n \times n}^t \quad (2.82)$$

2139 or

$$\mathbf{A}_{n \times n} = \sum_{k=1}^n \lambda_k \left(\mathbf{q}^{(k)} \right)_{n \times 1} \left((\mathbf{q}^{(k)})^t \right)_{1 \times n} \quad (2.83)$$

2140 This is a process of matrix decomposition by orthogonal matrices or by eigenvectors.
2141 Further, if all the eigenvalues are non-negative, then the matrix is said to be *positive*
2142 *semi-definite*, and the above formula can be written as

$$\mathbf{A}_{n \times n} = \sum_{k=1}^n \left(\mathbf{v}^{(k)} \right)_{n \times 1} \left((\mathbf{v}^{(k)})^t \right)_{1 \times n} \quad (2.84)$$

2143 where $\mathbf{v}^{(k)} = \sqrt{\lambda_k} \mathbf{q}^{(k)}$ ($k = 1, 2, \dots, n$). The sample covariance matrix in climate
2144 science satisfies the positive eigenvalue assumption, and will be discussed in more
2145 details in the next chapter. The last equation means that a positive semi-definite
2146 symmetric matrix $\mathbf{A}_{n \times n}$ can be decomposed into a sum of n outer products of
2147 eigenvectors.

2148 Matrix \mathbf{A}

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, \quad (2.85)$$

2149 is not positive semi-definite since its second eigenvalue is -1, but matrix \mathbf{C}

$$\mathbf{C} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad (2.86)$$

2150 is positive semi-definite since its eigenvalues are 3 and 1, which are positive. The
2151 matrix \mathbf{C} is actually positive definite. The eigenvectors are the same those of \mathbf{A} .
2152 Thus,

$$\mathbf{C} = \mathbf{Q} \mathbf{D}_c \mathbf{Q}^t, \quad (2.87)$$

2153 where \mathbf{Q} is given by Eq. (2.81), and \mathbf{D}_c is

$$\mathbf{D}_c = \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix}. \quad (2.88)$$

2154 The above can be verified by the following computer code.

```
2155 # Verify diagonalization and decomposition: R code
2156 C = matrix(c(2,1,1,2), nrow = 2)
2157 eigen(C)
```

```

2158  #$values
2159  #[1] 3 1
2160  #$vectors
2161  # [,1]      [,2]
2162  #[1,] 0.7071068 -0.7071068
2163  #[2,] 0.7071068  0.7071068
2164  Q = eigen(C)$vectors
2165  D = t(Q)%*%C%*%Q #Matrix diagonalization
2166  D
2167  #[1,]    3     0
2168  #[2,]    0     1
2169  Q%*%D%*%t(Q) #Matrix decomposition
2170  #[1,]    2     1
2171  #[2,]    1     2
2172  D[1,1]*Q[,1]%*%t(Q[,1]) + D[2,2]*Q[,2]%*%t(Q[,2])
2173  #[1,]    2     1
2174  #[2,]    1     2

```

```

# Verify diagonalization and decomposition: Python code
C = [[2,1],[1,2]]
valC, Q = np.linalg.eig(C)
print('eigenvalues of C = ', valC)
#eigenvalues of C = [3. 1.]
print('eigenvectors of C = ', Q)
#eigenvectors of C = [[ 0.70710678 -0.70710678]
# [ 0.70710678  0.70710678]]
D = Q.transpose(1,0).dot(C).dot(Q)
print('D = ', D)
#D = [[3. 0.]
# [0. 1.]]
# Matrix C is decomposed into three matrices: C = Q D Q'
Q.dot(D).dot(Q.transpose(1,0))
#array([[2., 1.],
#       [1., 2.]])
D[0][0]*np.outer(Q[:,0],Q.transpose()[:,0]) + \
D[1][1]*np.outer(Q[:,1],Q.transpose()[:,1])
#array([[ 1.,  2.],
#       [-2., -1.]]) #matrix C is recovered from vectors

```

2175

2176

2.5 Hadamard and Other Matrix Multiplications

2177

2178 In addition to the regular matrix products and sweeping by multiplication described
 2179 in Chapter 1, there still other matrix products, such as Hadamard product of two
 2180 matrices. This section presents a few more matrix operations.

2181 **2.5.1 Hadamard Product of Two Matrices of the Same
2182 Dimensions**

2184 Hadamard product of two matrices is the multiplication of the corresponding pair
2185 of elements of the two matrices of the same dimension:

$$A \circ B = [a_{ij}b_{ij}] \quad (2.89)$$

2186 The following is a computer code for an example of Hadamard product.

```
2187 #R Hadamard product of two matrices
2188 #install.packages('matrixcalc')
2189 library(matrixcalc)
2190 A = matrix( c( 1, 2, 3, 4 ), nrow=2, byrow=TRUE )
2191 B = matrix( c( 2, 4, 6, 8 ), nrow=2, byrow=TRUE )
2192 hadamard(A, B)
2193 # [,1] [,2]
2194 #[1,] 2 8
2195 #[2,] 18 32
```

2196 An application of the Hadamard product is to calculate the value of each item of
2197 a store, when the numbers of items for each good are stored as a matrix, and the
2198 corresponding unit prices as another. The Hadamard product of the two matrices
2199 yields a matrix of values for each kind of good.

2200 **2.5.2 Jordan Product of Two Matrices of the Same Dimensions**

2202 Jordan product of two square matrices of the same dimensions is defined as follows:

$$A \bullet B = \frac{1}{2}(AB + BA). \quad (2.90)$$

2203 A numerical example is as follows.

```
2204 #R Jordan product of A and B
2205 A = matrix( c( 1, 2, 3, 4 ), nrow=2, byrow=TRUE )
2206 B = matrix( c( 2, 1, 2, 1 ), nrow=2, byrow=TRUE )
2207 (A%*%B + B%*%A)/2
2208 # [,1] [,2]
2209 #[1,] 5.5 5.5
2210 #[2,] 9.5 7.5
```

2211 Since matrix multiplication is not commutative, $AB \neq BA$, Jordan product is
2212 the mean of AB and BA .

2213 **2.5.3 Commutator of Two Matrices of the Same Dimensions**

2215 The commutator of two matrices of the same dimensions is defined as follows:

$$[A, B] = AB - BA. \quad (2.91)$$

2216 A numerical example is below.

```

2217 #R commutator of A and B
2218 #R commutator of A and B
2219 A = matrix( c( 1, 2, 3, 4 ), nrow=2, byrow=TRUE )
2220 B = matrix( c( 2, 1, 2, 1 ), nrow=2, byrow=TRUE )
2221 A%*%B - B%*%A
2222 # [,1] [,2]
2223 #[1,] 1 -5
2224 #[2,] 9 -1
2225 #install.packages('psych')
2226 library(psych)
2227 tr(A%*%B - B%*%A) #tr for trace
2228 #[1] 0

```

As matrix multiplication is not commutative, $AB \neq BA$, the commutator measures the difference of AB minus BA . This operation has many applications in theoretical physics.

The trace of $[A, B]$ is zero. This can be easily proved as follows:

$$\text{tr}([A, B]) = \text{tr}(AB - BA) = \text{tr}(AB) - \text{tr}(BA) = \text{tr}(AB) - \text{tr}(AB) = 0. \quad (2.92)$$

Theorem 2.1 *Shoda theorem: If $\text{tr}(C)$ is zero, then C is a commutator matrix.*

The proof of this theorem can be made through constructing matrices A and B such that

$$C = AB - BA. \quad (2.93)$$

Given C , it is similar to another trace zero matrix D , i.e., there exists an orthogonal matrix S such as

$$D = SCS^{-1} = [d_{ij}]_{n \times n} \quad (2.94)$$

Then, it turns out that

$$A = SXS^{-1}, \quad B = SYS^{-1}, \quad (2.95)$$

where

$$X = \text{diag}(1, 2, \dots, n), \quad Y = [y_{ij}]_{n \times n} \quad (2.96)$$

with

$$y_{ij} = \begin{cases} (i-j)^{-1} d_{ij} & \text{else.} \\ 1 & \end{cases} \quad (2.97)$$

In fact,

$$D = [X, Y]. \quad (2.98)$$

2.5.4 Outer Product of Two Vectors and Two Matrices

2.5.4.1 Outer Product of Two Vectors

The outer product of two vectors a_m and $b_{p \times 1}$ is defined as

$$a \otimes b = a_{m \times 1} t(b)_{1 \times p} \quad (2.99)$$

and is an $m \times p$ matrix.

```

2247 #Outer product of two vectors
2248 a = 1:2
2249 b = 1:4
2250 a%o%b #outer product a_2-by-1 times t(b)_1-by-4
2251 # [1,] 1 2 3 4
2252 # [2,] 2 4 6 8
2253
2254
2255 #Outer product of A_mn and B_nq
2256 A = matrix(1:4, ncol = 2)
2257 B = matrix(1:6, ncol = 3)
2258 A%o%B
2259 dim(A%o%B)
2260 #[1] 2 2 2 3
2261
2262 #Outer product of A_mn and B_pq
2263 A = matrix(1:4, ncol = 2)
2264 B = matrix(1:9, ncol = 3)
2265 A%o%B
2266 dim(A%o%B)
2267 #[1] 2 2 2 3

```

This concept can be extended to the outer product of two matrices A and B , as shown in this numerical example.

2.5.4.2 Cross Product of Two Vectors in 3D

The cross product of two vectors $x = (x_1, x_2, x_3)$ and $y = (y_1, y_2, y_3)$ in 3-dimensional space is also a 3-dimensional vector and defined as follows

$$xy = (x_2y_3 - x_3y_2, x_3y_1 - x_1y_3, x_1y_2 - x_2y_1) \quad (2.100)$$

A numerical example is below.

```

2274 #R cross product and dot product
2275 library(pracma)
2276 x = 1:3
2277 y = 4:6
2278 cross(x, y)
2279 #[1] -3 6 -3
2280 dot(x, y)
2281 #[1] 32

```

This example also shows a dot product of two vectors. The dot product can be extended to n-dimensional space, but the cross product is limited to the 3-dimensional space. The cross product is used often in physics and engineering, particularly in mechanics and electromagnetism. The cross product of the force vector with an arm vector is a torque vector, which points to the direction of a screw when you drive it into a wall.

The dot product in 3D space also has many physics and engineering applications, such as the force's dot product with displacement vector is work.

2290 **2.5.4.3 Another Matrix Product Based on Outer Product of
2291 Vectors**

2292 There is another type of product of two matrices using the outer product of column
2293 and row vectors. We regard $A_{mn} = [a[, i]], i = 1, 2, \dots, n$ consisting of n column
2294 vectors, and $B_{n \times q} = [b[j,], j = 1, 2, \dots, n]$ consisting of n row vectors. The outer
2295 product is an mq matrix. defined as the

$$(AB)_{m \times q} = \left[\sum_{i=1}^n [a[, i]_{m \times 1} t(b[i,])_{1 \times q}] \right]_{m \times q}. \quad (2.101)$$

2296 **2.5.5 Kronecker Product of Two Matrices of the Same
2297 Dimensions**

2299 Kronecker product of two matrices A_{mn} and $B_{p \times q}$ is defined as the

$$A \diamond B = [a_{ij} B]_{mp \times nq}. \quad (2.102)$$

2300 Two numerical examples are as follows:

```
2301 #R Kronecker product of two matrices
2302 #Kronecker product
2303 library(fastmatrix)
2304 A <- diag(1:2)
2305 B <- matrix(1:4, ncol = 2)
2306 kronecker.prod(A, B)
2307 # [,1] [,2] [,3] [,4]
2308 #[1,] 1 3 0 0
2309 #[2,] 2 4 0 0
2310 #[3,] 0 0 2 6
2311 #[4,] 0 0 4 8
2312
2313 # an example with vectors
2314 ones <- rep(1, 2)
2315 y <- 1:4
2316 kronecker.prod(ones, t(y)) # 2-by-4 matrix
2317 # [,1] [,2] [,3] [,4]
2318 #[1,] 1 2 3 4
2319 #[2,] 1 2 3 4
```

2320 **Theorem 2.2** $(A \diamond B)^{-1} = A^{-1} \diamond B^{-1}$. **Proof:**

$$(A)(A^{-1} \diamond B^{-1}) = (AA^{-1}) \diamond (BB^{-1}) = I_n \diamond I_p = I_{np}. \quad (2.103)$$

2321 Here we have used the following equality

$$(A \diamond B)(C \diamond D) = (AC) \diamond (BD). \quad (2.104)$$

2322 This can be proved through direct matrix product computing.

2.6 Direct Sum of Two Matrices

The direct sum of matrices A and B is to stack the two matrices block-diagonally as illustrated in the following numerical example. The mathematical notation of the direct sum operation is often \oplus :

$$A \oplus B \quad (2.105)$$

```

2328 #R direct sum of two matrices
2329 A = matrix(1:4, ncol = 2)
2330 B = matrix(1:6, ncol = 3)
2331 A1 = rbind(A, matrix(rep(0, 4), ncol = 2))
2332 B1 = rbind(matrix(rep(0, 6), ncol = 3), B)
2333 C = cbind(A1, B1) #= direct sum of A and B
2334 C
2335 #      [,1] [,2] [,3] [,4] [,5]
2336 #[1,]    1     3     0     0     0
2337 #[2,]    2     4     0     0     0
2338 #[3,]    0     0     1     3     5
2339 #[4,]    0     0     2     4     6
2340
2341 #Express the direct sum by Kronecker products
2342 kronecker.prod(diag(1,0), A) + kronecker.prod(diag(0,1), B)
2343 #      [,1] [,2] [,3] [,4] [,5]
2344 #[1,]    1     3     0     0     0
2345 #[2,]    2     4     0     0     0
2346 #[3,]    0     0     1     3     5
2347 #[4,]    0     0     2     4     6

```

Theorem 2.3 *From Kronecker product to direct sum:*

$$A \oplus B = diag(1,0) \diamond A + diag(0,1) \diamond B. \quad (2.106)$$

The proof follows the definition of Kronecker product and has been verified in the numerical example.

2.7 Visualization of Eigenvalues and Eigenvectors for a Covariance Matrix

If $A_{n \times p}$ is a space-time data matrix, then

$$C_{n \times n} = \frac{1}{p} AA^t \quad (2.107)$$

is a covariance matrix. Here, n is the number of spatial locations, e.g., n clients of your health clinic; and p is apparently the number of time steps, e.g., the number of days under your supervision or treatment.

We wish to ask the following questions:

- 2359 (i) What are the typical differences between your clients?
 2360 (ii) What are the temporal variation patterns that show the beginning of the
 2361 effectiveness of your treatment?
 2362 The eigenvalues and eigenvectors of C may help answer Question (i). The SVD
 2363 in the next section may help answer Question (ii).
 2364 We use the following numerical example to illustrate the procedure of visualiza-
 2365 tion and interpretation of eigenvalues and eigenvectors for C .

2366 2.8 Singular Value Decomposition

2368 Previous section shows an eigenvector-eigenvalue decomposition of a symmetric
 2369 square matrix. A similar decomposition can be made for a rectangular matrix. The
 2370 decomposition using unit eigenvectors and eigenvalues for a general rectangular
 2371 matrix is called the *singular value decomposition* (SVD). Singular value is another
 2372 name of eigenvalue. Although the basic mathematical theory of SVD was devel-
 2373 oped almost 200 years ago (Stewart 1993), the modern algorithm of efficient SVD
 2374 computing was only developed by Gene H. Golub (1932-2007) and his colleagues in
 2375 the 1970s. Now, SVD has become an important data analysis tool for every field.
 2376 Climate science is not an exception. A space-time climate data matrix is often a
 2377 rectangular matrix, since the number of sites is not likely to be equal to the number
 2378 of temporal observations at those sites. We may denote a space-time climate data
 2379 matrix by $\mathbf{X}_{n \times m}$, where n denotes the number of sites, and m is the total number
 2380 of temporal observations. For $\mathbf{X}_{n \times m}$, we can also interpret n as the number of grid
 2381 boxes of a climate model output, and m as the number of time steps.

2382 2.8.1 SVD formula and a simple SVD example

2384 If $m \leq n$, then matrix $\mathbf{A}_{n \times m}$ has the following SVD decomposition

$$\mathbf{A}_{n \times m} = \mathbf{U}_{n \times m} \mathbf{D}_{m \times m} (\mathbf{V}^t)_{m \times m}. \quad (2.108)$$

2385 Here, \mathbf{U} may be interpreted as a spatial matrix, consisting of m spatial orthonormal
 2386 column vectors which are unit eigenvectors of $\mathbf{A}\mathbf{A}^t$; \mathbf{V} may be interpreted as a
 2387 temporal matrix, consisting of m temporal orthonormal column vectors which are
 2388 unit eigenvectors of $\mathbf{A}^t\mathbf{A}$; and \mathbf{D} is a diagonal matrix whose elements are the square
 2389 root of the eigenvalues of $\mathbf{A}\mathbf{A}^t$ and may be interpreted as standard deviations.
 2390 Figure 2.3 may help you understand this formula. The diagonal elements of \mathbf{D}
 2391 are called singular values, and the column vectors of \mathbf{U} and \mathbf{V} are called singular
 2392 vectors. Here, the word “singular” may be understood as special, or opposite to
 2393 “general”, or distinguished, or being out of ordinary.

2394 If $m \geq n$, then matrix $\mathbf{A}_{n \times m}$ has the following SVD decomposition

$$\mathbf{A}_{n \times m} = \mathbf{U}_{n \times n} \mathbf{D}_{n \times n} (\mathbf{V}^t)_{n \times m}. \quad (2.109)$$

SVD: $A = UDV^t$ when $n > m$ (top panel) or $n < m$ (bottom panel)

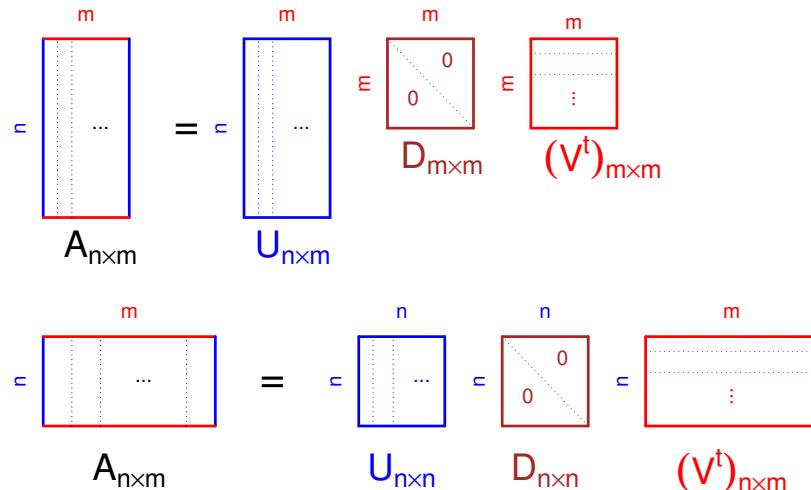


Fig. 2.3 Schematic diagrams of singular value decomposition (SVD): $A = UDV^t$.

2395 The following R code shows a simple SVD example of a 2×3 matrix.

```

2396 #SVD example for a 2-by-3 matrix: R code
2397 A=matrix(c(-1,1,0,2,-2,3),nrow=2)
2398 A #Show the 2-by-3 matrix
2399 # [,1] [,2] [,3]
2400 #[1,] -1 0 -2
2401 #[2,] 1 2 3
2402 svdA=svd(A) #Compute the SVD of A and put the results in svdA
2403 svdA #Show SVD results: d, U, and V
2404 round(svdA$d, digits=2) #Show only the singular values
2405 #[1] 4.22 1.09
2406 round(svdA$u, digits=2) #Show only matrix U
2407 # [,1] [,2]
2408 #[1,] -0.48 0.88
2409 #[2,] 0.88 0.48
2410 round(svdA$v, digits=2) #Show only matrix V
2411 # [,1] [,2]
2412 #[1,] 0.32 -0.37
2413 #[2,] 0.42 0.88
2414 #[3,] 0.85 -0.29
2415 sqrt(eigen(A%*%t(A))$values)
2416 #[1] 4.221571 1.085514

```

```

#SVD example for a 2-by-3 matrix: Python code
A = [[-1,0, -2],[1,2,3]]
UsvdA, DsvdA, VsvdA = np.linalg.svd(A)
print('Singular values= ', np.round(DsvdA,2))
#Singular values= [4.22 1.09]
print('Spatial singular vectors= ', np.round(UsvdA,2))
#Spatial singular vectors= [[-0.48 0.88]
# [ 0.88 0.48]]
print('Temporal singular vectors= ', np.round(VsvdA,2))
#Temporal singular vectors= [[ 0.32 0.42 0.85]
# [-0.37 0.88 -0.29]
# [-0.87 -0.22 0.44]]

B = np.array(A)
C = np.matmul(B, B.T) #B times B transpose
valC, vecC = np.linalg.eig(C)
np.sqrt(valC)
#[array([1.0855144 , 4.22157062])

```

2417

2418 The above computer code shows the following SVD results expressed in mathematical formulas below:

2420 (a) The vector form:

$$\begin{aligned}
 A &= \begin{bmatrix} -1 & 0 & -2 \\ 1 & 2 & 3 \end{bmatrix} \\
 &= 4.22 \begin{bmatrix} -0.48 \\ 0.88 \end{bmatrix} \times \begin{bmatrix} 0.32 & 0.42 & 0.85 \end{bmatrix} + \\
 &\quad 1.09 \begin{bmatrix} 0.88 \\ -0.48 \end{bmatrix} \times \begin{bmatrix} -0.37 & 0.88 & -0.29 \end{bmatrix}
 \end{aligned} \tag{2.110}$$

2421

2422 and

2423 (b) The matrix form:

$$A = \begin{bmatrix} -1 & 0 & -2 \\ 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} -0.48 & 0.88 \\ 0.88 & 0.48 \end{bmatrix} \begin{bmatrix} 4.22 & 0 \\ 0 & 1.09 \end{bmatrix} \begin{bmatrix} 0.32 & 0.42 & 0.85 \\ -0.37 & 0.88 & -0.29 \end{bmatrix} \tag{2.111}$$

2424 If we use only the first singular vectors to approximate A from the two triplets
2425 of singular vectors and singular values, the result is below.

```

#Data reconstruction by singular vectors: R code
round(svdA$d[1]*svdA$u[,1] %*% t(svdA$v[,1]),
       digits=1)
# [,1] [,2] [,3]
#[1,] -0.7 -0.8 -1.7
#[2,] 1.2  1.5  3.2

```

2432 It is quite close to A .

```

2433
2434 #Data reconstruction by singular vectors: Python code
2435 np.round(DsvdA[0]*np.outer(UsvdA[:,0], VsvdA[:,0]),1)
2436 #array([[-0.7, -0.8, -1.7],
2437 #        [ 1.2,  1.5,  3.2]])

```

2433

If we use both singular vectors to reconstruct A , then the reconstruction is exact without errors as expected:

```

2436 round(svdA$d[1]*svdA$u[,1]%^%t(svdA$v[,1]) +
2437   svdA$d[2]*svdA$u[,2]%^%t(svdA$v[,2]),
2438   digits =2)
2439 #[,1] [,2] [,3]
2440 #[1,] -1 0 -2
2441 #[2,] 1 2 3

```

2442

Figure 2.3 for the schematic diagram of SVD may be plotted by the following computer code.

```

2443 #R plot schematic diagram of SVD
2444 setwd('/Users/sshen/LinAlg')
2445 setEPS() #Plot the figure and save the file
2446 postscript("fig0503.eps", width = 11)
2447 par(mar=c(0,0,0,0))
2448 plot(200, axes = FALSE,
2449       xlab = "", ylab = "",
2450       xlim = c(-3,28), ylim = c(-3,16))
2451 text(13,15.5, cex=2.2,
2452       bquote("SVD:" ~ A == UDV^t ~ "when n > m or n < m"))
2453 #Space-time data matrix A when n>m
2454 segments(x0 = c(0,0,3,3),
2455            y0 = c(6,12,12,6) +1,
2456            x1 = c(0,3,3,0),
2457            y1 = c(12,12,6,6) +1,
2458            col = c('blue','red','blue','red')), lwd = 3)
2459 segments(x0 = c(0.5,1.0),
2460            y0 = c(6,6)+1,
2461            x1 = c(0.5,1.0),
2462            y1 = c(12,12)+1,
2463            lwd = 1.3, lty = 3)
2464 text(-.8, 9+1, 'n', srt=90, col ='blue', cex = 1.4)
2465 text(1.5, 12.8+1, 'm', col = 'red', cex = 1.4)
2466 text(2.0, 9+1, '...', cex = 1.4)
2467 text(2, 5+1, bquote(A[n %*% m]), cex = 2.5)
2468 text(5, 9+1, '=', cex = 3)
2469 #Spatial matrix U

```

```

2472 segments(x0 = c(7,7,10,10),
2473     y0 = c(6,12,12,6)+1,
2474     x1 = c(7,10,10,7),
2475     y1 = c(12,12,6,6)+1,
2476     col = c('blue','blue','blue','blue'), lwd =3)
2477 segments(x0 = c(7.5,8),
2478     y0 = c(6,6)+1,
2479     x1 = c(7.5,8),
2480     y1 = c(12,12)+1,
2481     lwd =1.3, lty = 3, col = 'blue')
2482 text(6.2, 9+1, 'n', srt=90, col = 'blue', cex = 1.4)
2483 text(8.5, 12.8+1, 'm', col = 'red', cex = 1.4)
2484 text(9, 9+1, '...', cex = 1.4, col='blue')
2485 text(8.7, 5.0+1, bquote(U[n%*%m]), cex = 2.5, col= 'blue')
2486 #Singular value diagonal matrix D
2487 segments(x0 = c(12,12,15,15),
2488     y0 = c(9,12,12,9)+1,
2489     x1 = c(12,15,15,12),
2490     y1 = c(12,12,9,9)+1,
2491     col = c('brown','brown','brown','brown'), lwd =3)
2492 segments(x0 = 12, y0 = 12+1, x1 = 15, y1 = 9+1, lty=3,
2493     col = c('brown'), lwd =1.3)#diagonal line
2494 text(11.2, 10.5+1, 'm', srt=90, col ='red', cex = 1.4)
2495 text(13.5, 12.8+1, 'm', col = 'red', cex = 1.4)
2496 text(14.1, 11.3+1, '0', col = 'brown', cex = 1.4)
2497 text(12.9, 10.0+1, '0', col = 'brown', cex = 1.4)
2498 text(13.9, 8.0+1, bquote(D[m%*%m]), cex = 2.5, col='brown')
2499 #Temporal matrix V
2500 segments(x0 = c(17,17,20,20),
2501     y0 = c(9,12,12,9)+1,
2502     x1 = c(17,20,20,17),
2503     y1 = c(12,12,9,9)+1,
2504     col = c('red','red','red','red'), lwd =3)
2505 segments(x0 = c(17,17),
2506     y0 = c(11.5,10.8)+1,
2507     x1 = c(20,20),
2508     y1 = c(11.5,10.8)+1,
2509     col = c('red','red'), lty=3, lwd =1.3)
2510 text(16.2, 10.5+1, 'm', srt=90, col ='red', cex = 1.4)
2511 text(18.5, 12.5+1, 'm', col = 'red', cex = 1.4)
2512 text(19.5, 8+1, bquote((V^t)[m%*%m]), cex = 2.5, col='red')
2513 text(18.5, 10+1, '...', col='red', srt=90, cex =1.4)
2514 # Space-time data matrix B when n < m
2515 segments(x0 = c(0,0,6,6),
2516     y0 = c(0,3,3,0),
2517     x1 = c(0,6,6,0),
2518     y1 = c(3,3,0,0),
2519     col = c('blue','red','blue','red'), lwd =3)
2520 segments(x0 = c(1,2,5),
2521     y0 = c(0,0,0),
2522     x1 = c(1,2,5),
2523     y1 = c(3,3,3),
2524     lwd =1.3, lty = 3)
2525 text(-0.8, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2526 text(3, 3.8, 'm', col = 'red', cex = 1.4)
2527 text(3.5, 1.5, '...', cex = 1.4)

```

```

2528 text(3, -1.5, bquote(A[n%*%m]), cex = 2.5)
2529 text(8, 1.5, '=', cex = 3)
2530 #Spatial matrix U
2531 segments(x0 = c(11,11,14,14),
2532             y0 = c(0,3,3,0),
2533             x1 = c(11,14,14,11),
2534             y1 = c(3,3,0,0),
2535             col = c('blue','blue','blue','blue'), lwd =3)
2536 segments(x0 = c(11.5,12.2),
2537             y0 = c(0,0),
2538             x1 = c(11.5,12.2),
2539             y1 = c(3,3),
2540             lwd =1.3, lty = 3, col = 'blue')
2541 text(10.2, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2542 text(12.5, 3.8, 'n', col = 'blue', cex = 1.4)
2543 text(13.2, 1.5, '...', cex = 1.4, col='blue')
2544 text(12.5, -1.5, bquote(U[n%*%n]), cex = 2.5, col= 'blue')
2545 #Singular value diagonal matrix D
2546 segments(x0 = c(16,16,19,19),
2547             y0 = c(0,3,3,0),
2548             x1 = c(16,19,19,16),
2549             y1 = c(3,3,0,0),
2550             col = c('brown','brown','brown','brown'), lwd =3)
2551 segments(x0 = 16, y0 = 3, x1 = 19, y1 = 0, lty=3,
2552             col = c('brown')), lwd =1.3)#diagonal line
2553 text(15.2, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2554 text(17.5, 3.8, 'n', col = 'blue', cex = 1.4)
2555 text(18.1, 2.3, '0', col = 'brown', cex = 1.4)
2556 text(16.9, 1.0, '0', col = 'brown', cex = 1.4)
2557 text(17.5, -1.5, bquote(D[n%*%n]), cex = 2.5, col='brown')
2558 #Temporal matrix V
2559 segments(x0 = c(21,21,27,27),
2560             y0 = c(0,3,3,0),
2561             x1 = c(21,27,27,21),
2562             y1 = c(3,3,0,0),
2563             col = c('red','red','red','red'),
2564             lwd =3)
2565 segments(x0 = c(21,21),
2566             y0 = c(2.5,1.8),
2567             x1 = c(27,27),
2568             y1 = c(2.5,1.8),
2569             col = c('red','red'), lty=3, lwd =1.3)
2570 text(20.2, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2571 text(24, 3.8, 'm', col = 'red', cex = 1.4)
2572 text(24, -1.5, bquote((V^t)[n%*%m]), cex = 2.5, col='red')
2573 text(24, 1, '...', col='red', srt=90, cex =1.4)
2574 dev.off()

```

```

2575 #Python plot Fig. 5.3: Schematic diagram of SVD
2576 import matplotlib.patches as patches
2577 import numpy as np
2578 import pylab as pl
2579 from matplotlib import collections as mc
2580 lines = [[(0, 7), (0, 13)], [(0, 13), (3, 13)],
2581         [(3, 13), (3, 7)], [(3, 7), (0, 7)]]
2582 c = np.array(['b', 'r', 'b', 'r'])
2583 lc = mc.LineCollection(lines, colors=c, linewidths=3)
2584 fig, ax = pl.subplots()
2585 ax.set_xlim([-3, 28])
2586 ax.set_ylim([-3, 16])
2587 ax.add_collection(lc)
2588 ax.margins(0.1)
2589 plt.plot([0.5, 0.5], [7, 13],
2590          linestyle='dotted', color = 'k')
2591 plt.plot([1, 1], [7, 13],
2592          linestyle='dotted', color = 'k')
2593 plt.text(13, 15.5,
2594           r'SVD: $A = UDV^T$ when $n > m$ or $n < m$',
2595           fontsize = 30)
2596 plt.text(-1.2, 10, 'n', color = 'blue',
2597           fontsize = 25, rotation=90)
2598 plt.text(1.1, 13.3, 'm', color = 'red',
2599           fontsize = 25, rotation=0)
2600 plt.text(0.0, 5.5, r'$A_{\{n \times m\}}$',
2601           fontsize = 35, rotation=0)
2602 plt.text(1.5, 10, '...', color = 'black',
2603           fontsize = 25, rotation=0)
2604 plt.axis('off')
2605 plt.show()

```

2575

2576 This Python code generates the top-left rectangular box in Fig. 5.3. The remaining code for other boxes is highly repetitive and can be found from the book website.

2579

2.9 SVD for the standardized sea level pressure data of Tahiti and Darwin

2580

2582 The Southern Oscillation Index (SOI) is an indicator for El Niño or La Niña. It is
2583 computed as the difference of sea level pressure (SLP) of Tahiti (17.75°S , 149.42°W)
2584 minus that of Darwin (12.46°S , 130.84°E). An SVD analysis of the SLP data can
2585 substantiate this calculation formula.

2586

The following shows the data matrix of the standardized SLP anomalies of Tahiti and Darwin from 2009 to 2015, and its SVD:

```

2588 #R SVD analysis for the weighted SOI from SLP data
2589 setwd("/Users/sshen/climmath")

```

```

2590 Pda<-read.table("data/PSTANDdarwin.txt", header=F)
2591 dim(Pda)
2592 #[1] 65 13 #Monthly Darwin data from 1951-2015
2593 pdaDec<-Pda[,13] #Darwin Dec standardized SLP anomalies data
2594 Pta<-read.table("data/PSTANDtahiti.txt", header=F)
2595 ptaDec=Pta[,13] #Tahiti Dec standardized SLP anomalies
2596 ptada1 = cbind(pdaDec, ptaDec) #space-time data matrix
2597
2598 #Space-time data format
2599 ptada = t(ptada1[59:65,]) #2009-2015 data
2600 colnames(ptada)<-2009:2015
2601 rownames(ptada)<-c("Darwin", "Tahiti")
2602 ptada #6 year of data for two stations
2603 # 2009 2010 2011 2012 2013 2014 2015
2604 #Darwin 0.5 -2.3 -2.2 0.3 0.3 0.1 -0.4
2605 #Tahiti -0.7 2.5 1.9 -0.7 0.4 -0.8 -1.3
2606 svdptd = svd(ptada) #SVD for the 2-by-6 matrix
2607 U=round(svdptd$u, digits=2)
2608 U
2609 #[1,] -0.66 0.75
2610 #[2,] 0.75 0.66
2611 D=round(diag(svdptd$d), digits=2)
2612 D
2613 #[1,] 4.7 0.00
2614 #[2,] 0.0 1.42
2615 V =round(svdptd$v, digits=2)
2616 t(V)
2617 #[1,] -0.18 0.72 0.61 -0.15 0.02 -0.14 -0.15
2618 #[2,] -0.06 -0.06 -0.28 -0.17 0.34 -0.32 -0.82

```

```

#Python SVD analysis of the Darwin and Tahiti SLP data
import os
os.chdir("/Users/sshen/LinAlg")
PDA = np.array(read_table("data/PSTANDdarwin.txt",
                           header = None, delimiter = "\s+"))
PTA = np.array(read_table("data/PSTANDtahiti.txt",
                           header = None, delimiter = "\s+"))
pdata = np.stack([PDA[58:65,12], PTA[58:65,12]], axis=0)
print('The Darwin and Tahiti SLP data 2009-2015 =', pdata)
#The Darwin and Tahiti Standardized SLP anomalies =
#[[ 0.5 -2.3 -2.2 0.3 0.3 0.1 -0.4]
# [-0.7 2.5 1.9 -0.7 0.4 -0.8 -1.3]]
u, d, v = np.linalg.svd(pdata)
print('Spatial singular vectors EOFs U =', np.round(u,2))
#Spatial singular vectors EOFs U = [[-0.66 0.75]
# [ 0.75 0.66]]
print('Diagonal matrix D =', np.round(np.diag(d),2))
#Diagonal matrix D = [[4.7 0.]
# [0. 1.42]]
print('Temporal singular vectors PCs V =', np.round(v,2))
#Temporal singular vectors PCs V=
#[[-0.18 0.72 0.61 -0.15 0.02 -0.14 -0.15]
# [-0.06 -0.06 -0.28 -0.17 0.34 -0.32 -0.82] ...]

```

2620 One can verify that

$$\mathbf{UDV}^t \quad (2.112)$$

2621 approximately recovers the original data matrix.

2622 The first column vector $(-0.66, 0.75)$ of the spatial pattern matrix \mathbf{U} may be
2623 interpreted to be associated with the Southern Oscillation Index (SOI), which puts
2624 a negative weight -0.66 on Darwin, and a positive weight 0.75 on Tahiti. The
2625 weighted sum is approximately equal to the difference of Tahiti's SLP minus that
2626 of Darwin, which the definition of SOI. The index measures large scale ENSO
2627 dynamics of the tropical Pacific (Trenberth 2020). The magnitude of the vector
2628 $(-0.66, 0.75)$ is approximately one, because U is a unitary matrix. The correspond-
2629 ing singular vector has a distinctly large value 0.72 in December 2010, which was
2630 a strong La Niña month. In this month, the Darwin has a strong negative SLP
2631 anomaly, while Tahiti has a strong positive SLP anomaly. This situation enhanced
2632 the easterly trade winds in the tropical Pacific and caused abnormally high precip-
2633 itation in Australia in the 2010-2011 La Niña period.

2634 The second column vector $(0.75, 0.66)$ of the spatial pattern matrix \mathbf{U} also has
2635 climate implications. The weighted sum with two positive equal weights measures
2636 the tropical Pacific dynamics of small scales (Trenberth 2020).

2637 2.10 Chapter summary

2638 A matrix can be regarded as a 2-dimensional $n \times m$ rectangular array of numbers
2639 or symbols, or as m column vectors, or as n row vectors. This may be interpreted as
2640 climate data at n locations with m time steps. Many mathematical properties of a
2641 matrix of climate data have climate interpretations. For example, SVD decomposes
2642 a space-time climate data matrix into an orthogonal spatial pattern matrix, an
2643 orthogonal temporal pattern matrix, and a diagonal “energy-level” matrix that
2644 measures the standard deviation of the temporal pattern:

$$\mathbf{A} = \mathbf{UDV}^t \quad (2.113)$$

2646 The column vectors of the spatial matrix \mathbf{U} are spatial singular vectors, also called
2647 EOFs, while those of the temporal matrix \mathbf{V} are temporal singular vectors, also
2648 called PCs. The first a few EOFs often have climate dynamic interpretations, such
2649 as El Niño Southern Oscillation (ENSO). If EOF1 corresponds to El Niño and
2650 shows some typical ENSO properties, such as the opposite signs of SLP anomalies
2651 of Darwin and Tahiti, then PC1 shows a temporal pattern, e.g., the extreme values
2652 of PC1 indicating both the occurrence time and the strength of El Niño. The
2653 diagonal elements of matrix \mathbf{D} are singular values, also known as eigenvalues.

2654 An eigenvector \mathbf{v} of a square matrix \mathbf{C} is such a special vector that \mathbf{C} 's action
2655 on \mathbf{v} does not change its orientation, i.e., \mathbf{Cv} is parallel to \mathbf{v} . This statement

2656 implies the existence of a scalar λ , called eigenvalue (also known as singular value
2657 or characteristic value), such that

$$\mathbf{C}\mathbf{v} = \lambda\mathbf{v}. \quad (2.114)$$

2658 We have also discussed the matrix method of solving a system of linear equations

$$\mathbf{Ax} = \mathbf{b}, \quad (2.115)$$

2659 linear independence of vectors, linear transform, and other basic matrix methods.
2660 These methods are useful for the chapters on covariance, EOFs, spectral analysis,
2661 regression analysis, and machine learning. You may focus on the computing meth-
2662 ods and computer code of the relevant methods. The mathematical proofs of this
2663 chapter, although helpful for exploring new mathematical methods, are not neces-
2664 sarily needed to read the other chapters of this book. If you are interested in an
2665 in-depth mathematical exploration of matrix theory, you may wish read the books
2666 by Horn and Johnson (1985) and Strang (2016).

References and Further Readings

- 2668 [1] Golub, G.H. and C. Reinsch, 1970: Singular value decomposition and least
 2669 squares solutions. *Numerische mathematik* 14, 403-420.

This seminal paper established an important method, known as the Golub-Reinsch algorithm, to compute the eigenvalues of a covariance matrix from a space-time matrix A without actually first computing the covariance matrix AA^t . This algorithm makes the SVD computation very efficient, which helps scientists consider SVD as a genuine linear algebra method, not a traditionally regarded statistical method based on a covariance matrix.

- 2670
 2671 [2] Horn, R. A., and C. R. Johnson, 1985: *Matrix Analysis*, Cambridge University
 2672 Press, 561pp.

This is a comprehensive book on matrix theory and is a good reference for a researcher in climate statistics. It assumes the knowledge of the first course of linear algebra.

- 2673
 2674 [3] Strang, G., 2016: *Introduction to Linear Algebra*. 5th edition, Wellesley-
 2675 Cambridge Press, Wellesley, MA 02482, 574pp.

Gilbert Strang (1934-) is an American mathematician and educator. His textbooks and pedagogy have been internationally influential. This text is one of the very few basic linear algebra books that includes excellent materials on SVD, probability, and statistics.

- 2676
 2677 [4] Stewart, G.W., 1993: On the early history of the singular value decomposition.
 2678 *SIAM Review* 35, 551-566.

This paper describes the contributions from five mathematicians in the period of 1814-1955 to the development of the basic SVD theory.

- 2679
 2680 [5] Trenberth, K., and National Center for Atmospheric Research Staff (Eds), 2020:
 2681 The Climate Data Guide: Southern Oscillation Indices: Signal, Noise and Tahiti-
 2682 ti/Darwin SLP (SOI). Retrieved from

2683 <https://climatedataguide.ucar.edu/climate-data/southern-oscillation-indices-signal-noise-and-tahitidarwin-slp-soi>

This site describes the optimal indices for large and small scale dynamics.

- 2685 [6] Tucker, A., 1993: The growing importance of linear algebra in undergraduate
2687 mathematics. *College Mathematics Journal* 24, 3-9.

This paper describes the historical development of linear algebra, such as the term “matrix” being coined by J.J. Sylvester in 1848, and pointed out that “tools of linear algebra find use in almost all academic fields and throughout modern society.” The use of linear algebra in the big data era is now even more popular.

2688

Exercises

- 2689 **2.1** Write a computer code to
2690 (a) Read the NOAAGlobalTemp data file, and
2693 (b) Generate a 4×8 space-time data matrix for the December mean surface
2694 air temperature anomaly data of four grid boxes and eight years.

2695 Hint: You may find the NOAA Global Surface Temperature (NOAAGlobal-
2696 Temp) dataset online. You can use either netCDF format or CSV format.

- 2697 **2.2** Write a computer code to find the inverse of the following matrix

```
#      [,1] [,2] [,3]
#[1,] 1.7 -0.7 1.3
#[2,] -1.6 -1.4 0.4
#[3,] -1.5 -0.3 0.6
```

- 2702 **2.3** Write a computer code to solve the following linear system of equations

$$\mathbf{Ax} = \mathbf{b}, \quad (2.116)$$

2703 where

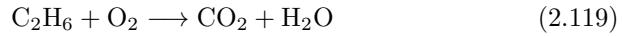
$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 0 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \quad (2.117)$$

- 2704 **2.4** The following equation

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad (2.118)$$

2705 has infinitely many solutions, and cannot be directly solved by a simple com-
 2706 puter command, such as `solve(A, b)`.
 2707 (a) Show that the three row vectors of the coefficient matrix are not linearly
 2708 independent.
 2709 (b) Because of the dependence, the linear system has only two independent
 2710 equations. Thus, reduce the linear system into two equations by treating x_3
 2711 as an arbitrary value while treating x_1 and x_2 as variables.
 2712 (c) Solve the two equations for x_1 and x_2 and express them in terms of x_3 .
 2713 The infinite possibilities of x_3 imply infinitely many solutions of the original
 2714 system.

2715 **2.5** Ethane is a gas similar to the greenhouse gas methane and can burn with
 2716 oxygen to form carbon dioxide and water:



2717 Given two ethane molecules, how many molecules of oxygen, carbon dioxide
 2718 and water should be in order for this chemical reaction equation to be bal-
 2719 anced?

2720 *Hint: Assume x, y and z molecules of oxygen, carbon dioxide and water, and
 2721 use the balance of the number of atoms of carbon, hydrogen and oxygen to
 2722 form linear equations. Solve the system of linear equations.*

2723 **2.6** Carry out the same procedure as the previous problem but for the burning of
 2724 methane CH_4 .

2725 **2.7** (a) Use matrix multiplication to show that the vector

$$\mathbf{u} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad (2.120)$$

2726 is not an eigenvector of the following matrix

$$\mathbf{A} = \begin{bmatrix} 0 & 4 \\ -2 & -7 \end{bmatrix} \quad (2.121)$$

2727 noindent (b) Find all the unit eigenvectors of matrix \mathbf{A} in Part (a).

2728 **2.8** Use hand-calculation to compute the matrix multiplication of \mathbf{UDV}^t where
 2729 the data of relevant matrices are given by the following R output:

```
2730 A=matrix(c(1,-1,1,1),nrow=2)
2731 A
2732 #[1,]    1    1
2733 #[2,]   -1    1
2734 svd(A)
2735 #\$d
2736 #[1] 1.414214 1.414214
2737 #\$u
2738 #[1,] -0.7071068 0.7071068
2739 #[2,]  0.7071068 0.7071068
2740 #\$v
2741 #[1,]   -1    0
2742 #[2,]    0    1
```

- 2743 **2.9** Use computer to find the matrices \mathbf{U} , \mathbf{D} and \mathbf{V} of the SVD for the following data matrix
2744

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 6 & 7 & 8 & 9 & 10 \\ 11 & 12 & 13 & 14 & 15 \end{bmatrix} \quad (2.122)$$

- 2745 **2.10** Use the first singular vectors to approximately reconstruct the data matrix
2746 \mathbf{A} using

$$B = d_1 \mathbf{u}_1 \mathbf{v}_1^t. \quad (2.123)$$

2747 Describe the goodness of the approximation using text, limited to 20 - 100
2748 words.

- 2749 **2.11** For the following data matrix

$$\mathbf{A} = \begin{bmatrix} 1.2 & -0.5 & 0.9 & -0.6 \\ 1.0 & -0.7 & -0.4 & 0.9 \\ -0.2 & 1.1 & 1.6 & -0.4 \end{bmatrix}, \quad (2.124)$$

- 2750 (a) Use computer to find the eigenvectors and eigenvalues of matrix $\mathbf{A}\mathbf{A}^t$.
2751 (b) Use computer to find the eigenvectors and eigenvalues of matrix $\mathbf{A}^t\mathbf{A}$.
2752 (c) Use computer to calculate SVD of \mathbf{A} .
2753 (d) Compare the singular vectors and singular values in Step (c) with the
2754 eigenvalues and eigenvectors computed in Steps (a) and (b). Use text to discuss
2755 your comparison.
2756 (e) Use computer to verify that the column vectors of \mathbf{U} and \mathbf{V} in the SVD
2757 of Step (c) are orthonormal to each other.

- 2758 **2.12** Make an SVD analysis of the December standardized anomalies data of sea
2759 level pressure (SLP) at Darwin and Tahiti from 1961 to 2010. You can find
2760 the data from the Internet or from the website of this book.
2761 (a) Write a computer code to organize the data into a 2×50 space-time data
2762 matrix.
2763 (b) Make the SVD calculation for this space-time matrix.
2764 (c) Plot the first singular vector in \mathbf{V} against time 1961 to 2010 as a time
2765 series curve, which is called the first principal component, denoted by PC1.
2766 (d) Interpret the first singular vector in \mathbf{U} , which is called the first empirical
2767 orthogonal function (EOF1), as weights of Darwin and Tahiti stations.
2768 (e) Check the historical El Niño events between 1961 and 2010 from the
2769 Internet, and interpret the extreme values of PC1.

- 2770 **2.13** Plot PC2 against time from the SVD analysis in the previous problem. Discuss
2771 the singular values λ_1 and λ_2 and interpret PC2, in reference to the description
2772 entitled “The Climate Data Guide: Southern Oscillation Indices: Signal, Noise
2773 and Tahiti/Darwin SLP (SOI)” by Kevin Trenberth (2020).
2774 **2.14** Make an SVD analysis similar to the previous two problems for the January
2775 standardized SLP anomalies data at Darwin and Tahiti from 1961 to 2010.

- 2776 **2.15** Make an SVD analysis similar to the previous problem for the monthly stan-
 2777 dardized SLP anomalies data at Darwin and Tahiti from 1961 to 2010. This
 2778 problem includes anomalies for every month. The space-time data is a 2×600
 2779 matrix.
- 2780 **2.16** For the observed data of the monthly surface air temperature at five stations
 2781 of your interest from January 1961 to December 2010, form a space-time data
 2782 matrix, compute the SVD of this matrix, and interpret your results from the
 2783 perspective of climate science. Please plot PC1, PC2, and PC3 against time.
 2784 You may find your data from the Internet, such as the NOAA Climate Data
 2785 Online website <https://www.ncdc.noaa.gov/cdo-web>
- 2786 **2.17** Do the same analysis as the previous problem, but for the monthly precipita-
 2787 tion data at the same stations in the same time period.
- 2788 **2.18** For an Reanalysis dataset, make an SVD analysis similar to the previous
 2789 problem for the monthly surface temperature data over 10 grid boxes of your
 2790 choice in the time period of 1961-2010. The space-time data is a 10×600
 2791 matrix. You can use your preferred Reanalysis dataset and download the data
 2792 from the Internet, such as NCEP/NCAR Reanalysis, and ECMWF ERA.
- 2793 **2.19** Let $\mathbf{C} = \mathbf{AA}^t$ and \mathbf{A} is any real-valued rectangular matrix. Show that
 2794 (a) the eigenvalues of \mathbf{C} are non-negative;
 2795 (b) if $\mathbf{v} = \mathbf{A}^t \mathbf{u}$, then \mathbf{v} is a eigenvector of $\mathbf{C}^t = \mathbf{A}^t \mathbf{A}$.
- 2796 **2.20** Given that

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad (2.125)$$

2797 write down the second order polynomial corresponding to the following matrix
 2798 expression

$$P(x_1, x_2) = \mathbf{x}^t \mathbf{A} \mathbf{A}^t \mathbf{x}. \quad (2.126)$$

- 2799 This is known as the quadratic form of the matrix $\mathbf{A} \mathbf{A}^t$.
- 2800 **2.21** If \mathbf{x} is a unit vector, use calculus to find the maximum value of $P(x_1, x_2)$ in
 2801 the previous problem. How is your solution related to the eigenvalue of $\mathbf{A} \mathbf{A}^t$.

Matrix Applications to Machine Learning

2804
2805 Machine learning (ML) is a branch of science that uses data and algorithms to
2806 mimic how human beings learn. The accuracy of the ML results can be gradually
2807 improved based on new training data and algorithm update. For example, a baby
2808 learns how to pick an orange from a fruit plate containing apples, bananas and
2809 oranges. Another baby learns how to sort out different kinds of fruits from a basket
2810 into three categories without naming the fruits. Then, how does ML work? It is
2811 basically a decision process for clustering, classification, or prediction, based on
2812 the input data, decision criteria, and algorithms. It does not stop here. It further
2813 validates the decision results and quantifies errors. The errors and the updated data
2814 will help update the algorithms and improve the results.

2815 ML has recently become a very popular method in climate science due to the
2816 availability of powerful and convenient resources of computing. It has been used to
2817 predict weather and climate, and to develop climate models. This chapter is a brief
2818 introduction of ML and provides basic ideas and examples. Our materials will help
2819 readers understand and improve the more complex ML algorithms used in climate
2820 science, so that they can go a step beyond only applying the ML software packages
2821 as a black box. We also provide R and Python codes for some basic ML algorithms,
2822 such as K-means for clustering, support vector machine for the maximum separation
2823 of sets, random forest of decision trees for classification and regression, and neural
2824 network training and predictions.

2825 Artificial intelligence (AI) allows computers to automatically learn from past
2826 data without human programming, which enables a machine to learn and to have
2827 intelligence. Machine learning is a subset of AI. Our chapter here focuses on ML,
2828 not the general AI.

2829 3.1 K-means clustering

2830
2831 A few toddlers at a daycare center may learn how to grab a few candies near
2832 themselves. It can be a point of fighting for a candy at a location that is not
2833 obviously closer to one than another. K-means method can help divide the candies
2834 among the toddlers in a fair way.

2835 Based on the historical weather data over a country, can ML decide the climate
2836 regimes for the country? Can ML determine the ecoregions of a country? Can ML

2837 define the regimes of wild fire over a region? The K-means clustering method can
 2838 be useful to answering these questions.

2839 2840 3.1.1 K-means setup and trivial examples

2841 The aim of the K-means clustering is to divide N points into K clusters so that
 2842 the total within cluster sum of squares (tWCSS) is minimized. Here we use 2D
 2843 data points to describe tWCSS and the K-means algorithm, although the K-means
 2844 method can be formulated in higher dimensions. We regard the data $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$
 2845 as the 2D coordinates of N points. For example, we treat $\mathbf{x}_1 = (1.2, -7.3)$ as obser-
 2846 vational data. Assume that these points can be divided into K clusters (C_1, C_2, \dots, C_K) ,
 2847 where K is subjectively given by you, the user who wishes to divide the N points
 2848 into K clusters by the K-means method. Let $\mathbf{x} \in C_i$, i.e., points within cluster C_i .
 2849 The number of points in cluster C_i is unknown and is to be determined by the
 2850 K-means algorithm. The total WCSS is defined by the following formula

$$\text{tWCSS} = \sum_{i=1}^K \left(\sum_{\mathbf{x} \in C_i} |\mathbf{x} - \boldsymbol{\mu}_i|^2 \right) \quad (3.1)$$

2851 where

$$\boldsymbol{\mu}_i = \frac{1}{K_i} \sum_{\mathbf{x} \in C_i} \mathbf{x} \quad (3.2)$$

2852 is the mean of the points within the cluster C_i , as K_i is the number of points in C_i .
 2853 Thus, we have K means, the name of the K-means method. That in the parentheses
 2854 is the within cluster sum of squares (WCSS)

$$\text{WCSS}_i = \left(\sum_{\mathbf{x} \in C_i} |\mathbf{x} - \boldsymbol{\mu}_i|^2 \right). \quad (3.3)$$

2855 This is defined for each cluster, and tWCSS is the sum of these WCSS_i for all the
 2856 K clusters.

2857 Some publications use WCSS or Tot WCSS, instead of tWCSS, to express the
 2858 right hand side of Eq. (3.1). Be careful when reading literatures concerning the
 2859 WCSS definition. Some computer software for K-means may even use a different
 2860 tWCSS definition.

2861 Our K-means computing algorithm is to minimize the tWCSS by optimally or-
 2862 ganizing the N points into K clusters. This algorithm can assign each data point
 2863 to a cluster. If we regard $\boldsymbol{\mu}_i$ as the centroid or center of cluster C_i , we may say that
 2864 the points in cluster C_i may have some kind of similarity, such as similar climate or
 2865 similar ecological characteristics, based on certain criteria. Below we use two simple
 2866 cases ($N = 2$ and $N = 3$) to illustrate the K-means method and its solutions.

2867 (i) Case $N = 2$:

2868

2869 When $N = 2$, if we assume $K = 2$, then

$$\mu_1 = \mathbf{x}_1, \quad \mu_2 = \mathbf{x}_2 \quad (3.4)$$

2870 and

$$\text{tWCSS} = \sum_{i=1}^2 \sum_{\mathbf{x} \in C_i} |\mathbf{x} - \mu_i|^2 = 0. \quad (3.5)$$

2871 This is the global minimum for $K = 2$, and

$$\mathbf{x}_1 \in C_1, \quad \mathbf{x}_2 \in C_2, \quad \mu_1 = \mathbf{x}_1, \quad \mu_2 = \mathbf{x}_2 \quad (3.6)$$

2872 is the unique solution. Of course, this is a trivial solution and bears no meaning.

2873 When $N = 2$, if we assume $K = 1$, then

$$\mu_1 = \frac{\mathbf{x}_1 + \mathbf{x}_2}{2} \quad (3.7)$$

2874 and

$$\text{tWCSS} = \sum_{i=1}^1 \sum_{\mathbf{x} \in C_i} |\mathbf{x} - \mu_i|^2 = \frac{1}{4} |\mathbf{x}_1 - \mathbf{x}_2|^2. \quad (3.8)$$

2875 This is the unique minimum of tWCSS, and the solution is also trivial.

2876 (ii) Case $N = 3$:

2877 When $N = 3$, if we assume $K = 2$, then there are three different cases of clustering

$$C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}, \quad C_2 = \{\mathbf{x}_1, \mathbf{x}_3\}, \quad C_3 = \{\mathbf{x}_2, \mathbf{x}_3\}. \quad (3.9)$$

2878 Here, the case

$$C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}, \quad C_2 = \{\mathbf{x}_3\} \quad (3.10)$$

2879 and

$$C_2 = \{\mathbf{x}_1, \mathbf{x}_2\}, \quad C_1 = \{\mathbf{x}_3\} \quad (3.11)$$

2880 are considered the same, since their difference is only a matter which cluster is called C_1 or C_2 .

2881 For case $C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}$, $C_2 = \{\mathbf{x}_3\}$, we have

$$\mu_1 = \frac{1}{2}(\mathbf{x}_1 + \mathbf{x}_2), \quad \mu_2 = \mathbf{x}_3, \quad (3.12)$$

2882 and

$$\text{tWCSS} = \frac{1}{4} |\mathbf{x}_1 - \mathbf{x}_2|^2. \quad (3.13)$$

2883 The tWCSS for the other two cases of clustering ($C_1 = \{\mathbf{x}_1, \mathbf{x}_3\}$, $C_2 = \{\mathbf{x}_2, \mathbf{x}_3\}$) can be computed in a similar way. Then, the smallest tWCSS from the three determines the final solution, which is one of the three cases of clustering. A numerical example is given below.

2889

2890

Example 3.1 Find the two K-means clusters for the following three points

$$P_1(1, 1), \quad P_2(2, 1), \quad P_3(3, 3.5). \quad (3.14)$$

2891 This is a K-means problem with $N = 3$ and $K = 2$. The K-means clusters can be
 2892 found by a computer command, e.g., `kmeans(mydata, 2)` in R. Figure 3.1 shows
 2893 the two K-means clusters with their centers at

$$C_1(1.5, 1), \quad C_2(3, 3.5) \quad (3.15)$$

2894 and

$$\text{tWCSS} = 0.5. \quad (3.16)$$

2895 Points P_1 and P_2 , indicated by red dots, are assigned to cluster C_1 whose center is
 2896 marked by a red cross. Point P_3 , indicated by the sky blue dot, is assigned to cluster
 2897 C_2 , whose center is marked with a blue cross and obviously overlaps with P_3 . This
 2898 result agrees with our intuition since P_1 and P_2 are together. The three points can
 2899 have two other possibilities of combination $C_1 = \{P_1, P_3\}$ and $C_1 = \{P_2, P_3\}$. The
 2900 case $C_1 = \{P_1, P_3\}$ leads to $\text{tWCSS} = 5.125$, and $C_1 = \{P_2, P_3\}$ to $\text{tWCSS} = 3.625$.
 2901 Both tWCSS are greater than 0.5. Thus, these two combinations should not be a
 2902 solution of the K-means clustering. These numerical results may be produced by
 2903 the following computer codes.

```

2904 #tWCSS calculation for N = 3 and K = 2
2905 mydata <- matrix(c(1, 1, 2, 1, 3, 3.5),
2906                     nrow = N, byrow = TRUE)
2907 x1 = mydata[1, ]
2908 x2 = mydata[2, ]
2909 x3 = mydata[3, ]
2910
2911 #Case C1 = (P1, P2)
2912 c1 = (mydata[1, ] + mydata[2, ])/2
2913 c2 = mydata[3, ]
2914 tWCSS = norm(x1 - c1, type = '2')^2 +
2915   norm(x2 - c1, type = '2')^2 +
2916   norm(x3 - c2, type = '2')^2
2917 tWCSS
2918 #[1] 0.5
2919
2920 #Case C1 = (P1, P3)
2921 c1 = (mydata[1, ] + mydata[3, ])/2
2922 c2 = mydata[2, ]
2923 norm(x1 - c1, type = '2')^2 +
2924   norm(x3 - c1, type = '2')^2 +
2925   norm(x2 - c2, type = '2')^2
2926 #[1] 5.125
2927
2928 #Case C1 = (P2, P3)
2929 c1 = (mydata[2, ] + mydata[3, ])/2
2930 c2 = mydata[1, ]
2931 norm(x2 - c1, type = '2')^2 +
2932   norm(x3 - c1, type = '2')^2 +

```

```
2933     norm(x1 - c2, type = '2')^2
2934 #[1] 3.625
2935
2936 #The case C1 = (P1, P2) can be quickly found by
2937 kmeans(mydata, 2)
2938 #Clustering vector:
2939 #[1] 1 1 2 #points P1, P2 in C1
```

```

#Python code: K-means computing for N = 3 and K = 2
N = 3
# create a 3 x 2 matrix of data
mydata = np.array([1,1,2,1,3,3.5]).reshape(3,2)

# first row of data
x1 = mydata[0,:]
# second row of data
x2 = mydata[1,:]
# third row of data
x3 = mydata[2,:]

# case C1 = (P1, P2)
c1 = (x1 + x2)/2
c2 = mydata[2,:]
tWCSS = np.linalg.norm(x1 - c1)**2 +
np.linalg.norm(x2 - c1)**2 + \
np.linalg.norm(x3 - c2)**2
print(tWCSS)
#0.5

# case C1 = (P1, P3)
c1 = (x1 + x3)/2
c2 = mydata[1,:]
print(np.linalg.norm(x1 - c1)**2 +
np.linalg.norm(x3 - c1)**2 + \
np.linalg.norm(x2 - c2)**2)

# case C1 = (P2, P3)
c1 = (x2 + x3)/2
c2 = mydata[0,:]
print(np.linalg.norm(x2 - c1)**2 +
np.linalg.norm(x3 - c1)**2 + \
np.linalg.norm(x1 - c2)**2)
#5.125

# case C1 = (P1, P2) can be quickly found by
kmeans = KMeans(n_clusters=2).fit(mydata)

print(kmeans.labels_)
#[0 0 1] #points P1, P2 in C1
#because Python index begins with 0 while R starts from 1

# number of data points
N = 3
# assume K clusters
K = 2
# create a 3 x 2 matrix of data
mydata = np.array([1, 1, 2, 1, 3, 3.5]).reshape(3,2)

Kclusters = KMeans(n_clusters=K).fit(mydata)
# cluster means
print(Kclusters.cluster_centers_, "\n")
#[[1.5 1. ]
#[3. 3.5]]
# sum of squares by cluster
print(Kclusters.inertia_)
#0.5

```

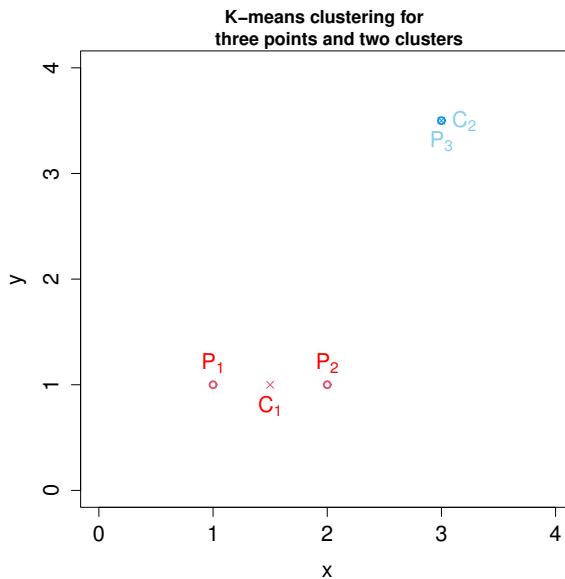


Fig. 3.1 The K-means clustering results for three points ($N = 3$) and two clusters ($K = 2$).

2941 Of course, this problem is trivial and can even be solved analytically by hand.
 2942 The computer solution of the K-means clustering is necessary when more points
 2943 are involved and when it is not obvious for a point to belong to a specific cluster.
 2944

2945 Figure 3.1 may be generated by the following computer code.

```

2946 # R plot Fig. 9.1: K-means for N = 3 and K = 2
2947 setwd("~/climstats")
2948 N = 3 #The number of data points
2949 K = 2 #Assume K clusters
2950 mydata = matrix(c(1, 1, 2, 1, 3, 3.5),
2951                   nrow = N, byrow = TRUE)
2952 Kclusters = kmeans(mydata, K)
2953 Kclusters #gives the K-means results,
2954 #e.g., cluster centers and WCSS
2955 #Cluster means:
2956 #[,1] [,2]
2957 #1 1.5 1.0
2958 #2 3.0 3.5
2959 #Within cluster sum of squares by cluster:
2960 #[1] 0.5 0.0
2961 Kclusters$centers
2962 par(mar = c(4,4,2.5,0.5))
2963 plot(mydata[,1], mydata[,2], lwd = 2,
2964       xlim = c(0, 4), ylim = c(0, 4),
2965       xlab = 'x', ylab = 'y', col = c(2, 2, 4),
2966       main = 'K-means clustering for
2967       three points and two clusters',

```

```

2968      cex.lab = 1.4, cex.axis = 1.4)
2969 points(Kclusters$centers[,1], Kclusters$centers[,2],
2970           col = c(2, 4), pch = 4)
2971 text(1.5, 0.8, bquote(C[1]), col = 'red', cex = 1.4)
2972 text(3.2, 3.5, bquote(C[2]), col = 'skyblue', cex = 1.4)
2973 text(1, 1.2, bquote(P[1]), cex = 1.4, col = 'red')
2974 text(2, 1.2, bquote(P[2]), cex = 1.4, col = 'red')
2975 text(3, 3.3, bquote(P[3]), cex = 1.4, col = 'skyblue')

```

```

#Python plot Fig. 9.1: SVM clustering for N=3 and K=2
#starting from setting up the plot
plt.figure(figsize = (8,8))
plt.ylim([0,4])
plt.yticks([0,1,2,3,4])
plt.xlim([0,4])
plt.xticks([0,1,2,3,4])
plt.title("K-means clustering for \n\u2022\u2022\u2022 three points and two clusters", pad = 10)
plt.xlabel("x", labelpad = 10)
plt.ylabel("y", labelpad = 10)

# plot P1-P3
plt.scatter(mydata[:,0], mydata[:,1],
            color = ('r', 'r', 'dodgerblue'),
            facecolors='none', s = 80)
# plot C1 and C2
plt.scatter((Kclusters.cluster_centers_[0][0],
             Kclusters.cluster_centers_[1][0]),
            (Kclusters.cluster_centers_[0][1],
             Kclusters.cluster_centers_[1][1]),
            marker = 'X', color = ('dodgerblue', 'r'))

# add labels
plt.text(1.43, 0.8, "$C_1$", color = 'r', size = 14)
plt.text(3.1, 3.45, "$C_2$", color = 'dodgerblue', size = 14)
plt.text(0.95, 1.1, "$P_1$", color = 'r', size = 14)
plt.text(1.95, 1.1, "$P_2$", color = 'r', size = 14)
plt.text(2.95, 3.3, "$P_3$", color = 'dodgerblue', size = 14)

plt.show()

```

2976

2977
2978

3.1.2 The number of exhaustive K clusters from N points

2979 The search for the minimum tWCSS in the previous sub-section for three points is
2980 through an algorithm of exhaustive search. This brute-force algorithm enumerates
2981 all possible cases. This algorithm can only work for a small N . When N is large, the
2982 number of cases will become too large to be calculated by a computer. For example,
2983 if $N = 100$ and $K = 10$, then the number of different clusters is approximately equal
2984 to

$$C(99, 9) = 1.7310 \times 10^{12}. \quad (3.17)$$

2985 It is impossible for a computer to exhaust all these cases to compute tWCSS for
 2986 each case. The computer code for the combination is

2987

```
2988 choose(99, 9) #R code to choose 9 out of 99
2989 #[1] 1.731031e+12 for N = 100, K =10
```

2990 The number of cluster cases can be computed following the so-called stars-and-
 2991 bars theorem, also known as the sticks-and-stones theorem, in counting. The stars-
 2992 and-bars theorem states that when using $K - 1$ bars or sticks to separate a line of
 2993 N stars, what is the number of possible cases of separation? Or put it in another
 2994 way. N number of candies is divided among K children. How many ways to divide
 2995 the candies, assuming that each child will receive at least one candy? The answer
 2996 is

$$C(N - 1, K - 1), \quad (3.18)$$

2997 because you can line up all the N stars and $K - 1$ bars among them to form K
 2998 subsets of stars. You first put one star in each subset. After this, you have $N - K$
 2999 stars to divide among the K subsets, i.e., K clusters. The $N - K$ stars plus $K - 1$
 3000 bars form $(N - K) + (K - 1) = N - 1$ objects that include both stars and bars.
 3001 Then the problem becomes choosing $K - 1$ objects from $N - 1$ objects, which has
 3002 $C(N - 1, K - 1)$ cases.

3003
3004

3.1.3 A K-means algorithm

3005 The previous sub-section has demonstrated that the brute-force exhaustive search
 3006 is not an algorithm choice to find the best K clusters from N points. We introduce
 3007 the K-means algorithm in this sub-section.

3008 From the above sub-section, we may have already developed a feeling that the
 3009 K-means clustering method may roughly be divided into the following two steps:

- 3010 • Assignment Step: Assume a number K and randomly generate K points $\mu_i^{(1)}$ ($i =$
 3011 $1, 2, \dots, K$) as the initial K centroids for K clusters. Without replacement, for
 3012 each point \mathbf{x} , compute the distance

$$d_i = |\mathbf{x} - \mu_i^{(1)}| \quad i = 1, 2, \dots, K. \quad (3.19)$$

3013 Assign the point \mathbf{x} to the cluster with the smallest distance. Then, assign the
 3014 next point to a cluster until every point is assigned. Thus, the initial K clusters
 3015 are formed: $C_i^{(1)}$ and each $C_i^{(1)}$ containing $K_i^{(1)}$ data points ($i = 1, 2, \dots, K$).

- 3016 • Update Step: Compute the centroids of the initial K clusters

$$\mu_i^{(2)} = \frac{1}{K_i^{(1)}} \sum_{\mathbf{x} \in C_i^{(1)}} \mathbf{x} \quad (3.20)$$

3017 and compute tWCSS⁽²⁾

$$\text{tWCSS}^{(2)} = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i^{(2)}} |\mathbf{x} - \boldsymbol{\mu}_i^{(2)}|^2 \quad (3.21)$$

3018 These $\boldsymbol{\mu}_i^{(2)}$ ($i = 1, 2, \dots, K$) are regarded as the updated centroids. We then
 3019 repeat the Assignment Step to assign each point to a cluster again according
 3020 to the rule of the nearest distance. Next, we compute $\boldsymbol{\mu}_i^{(3)}$ ($i = 1, 2, \dots, K$),
 3021 and tWCSS⁽³⁾

3022 The K-means algorithm continues the above iterative steps of Assignment and
 3023 Update until the centroids $\boldsymbol{\mu}_i^{(n)}$ ($i = 1, 2, \dots, K$) do not change from those in the
 3024 previous assignments $\boldsymbol{\mu}_i^{(n-1)}$ ($i = 1, 2, \dots, K$). When the assignment centroids do
 3025 not change, we say that the K-means iterative sequence has converged. It can be
 3026 rigorously proved that this algorithm always converges.

3027 R and Python have K-means functions. Many software packages are freely available
 3028 for some modified K-means algorithms using different distance formulas. Climate
 3029 scientists often just use these functions or packages without actually writing
 3030 an original computer code for a specific K-means algorithm. For example, the R
 3031 calculation in the previous sub-section used the command `kmeans(mydata, K)` to
 3032 calculate the cluster results for three points shown in Fig. 3.1.

3033 The K-means algorithm appears very simple, but it was invented only in the
 3034 1950s. Although the above K-means iterative sequence always converges, the final
 3035 result may depend on the initial assignment of the centroid: $C_i^{(1)}$ ($i = 1, 2, \dots, K$),
 3036 hence each run of the R K-means command `kmeans(mydata, K)` may not always
 3037 lead to the same solution. Further, a K-means clustering result may not be a global
 3038 minimum for tWCSS.

3039 Thus, after we have obtained our K-means results, we should try to explain the
 3040 results and see if they are reasonable. We may also run the code with different
 3041 K values and check which value K is the most reasonable. The tWCSS score is a
 3042 function of K . We may plot the function tWCSS(K) and check the shape of the
 3043 curve. Usually, tWCSS(K) is a decreasing function, and further,

$$\Delta_K = \text{tWCSS}(K) - \text{tWCSS}(K + 1) > 0, \quad K = 1, 2, \dots \quad (3.22)$$

3044 decreases fast for $K \leq K_e$, and then the decrease suddenly slows down or even
 3045 increases at $K = K_e$. In this sense, the curve tWCSS(K) looks like having an elbow
 3046 at K_e . We would choose this K_e as the optimal K . See the elbow in the example
 3047 of the next sub-section.

3048 Different applications may need different optimization criteria for the selection
 3049 of the best number of clusters. When you work on your own data, you may try
 3050 different criteria to choose the best number of clusters. For example, the Silhouette
 3051 method uses the Silhouette score to measure the similarity level of a point with its
 3052 own cluster compared to other clusters.

3053 A computer code for the K-means clustering may not always obtain the minimum

3054 tWCSS against all partitions. Sometimes, a computer yields a “local” minimum for
 3055 tWCSS(K).

3056 In short, when using the K-means method to divide given points into clusters,
 3057 you should run the K-means code multiple times, check the relevant solution results.
 3058 Although most K-means code yields a unique solution, multiple solutions sometimes
 3059 appear, then we need to examine which solution is the most reasonable for our
 3060 purposes.

3061 3.1.4 K-means clustering for the daily Miami weather data

3063 In this sub-section, we present an example using the K-means method to cluster
 3064 the observed daily weather data at the Miami International Airport, the United
 3065 States.

3066 3.1.4.1 K-means clustering for $N = 365, K = 2$

3067 Figure 3.2(a) shows a scatter plot of the daily minimum surface air temperature
 3068 T_{\min} [in $^{\circ}\text{C}$] and the direction of the fastest 2-minute wind in a day WDF2 [in
 3069 degrees] of the Miami International Airport. The wind blowing from the north is
 3070 zero degree, from the east is 90 degrees, and from the west is 270 degrees. The
 3071 scatter plot seems to support two clusters or weather regimes: The lower regime
 3072 showing the wind direction between 0 and 180 degrees, meaning the wind from the
 3073 east, northeast, or southeast; and another with the wind direction between 180
 3074 degrees and 360 degrees, meaning that the wind came from the west, northwest, or
 3075 southwest;

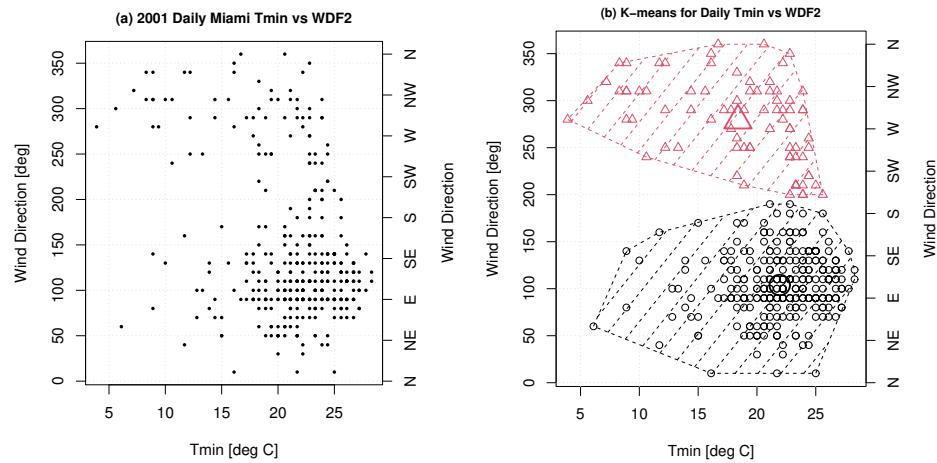


Fig. 3.2 (a) Scatter plot of the daily T_{\min} vs WDF2 for the Miami International Airport in 2001. (b) The K-means clusters of the the daily T_{\min} vs WDF2 data points when assuming $K = 2$.

3076 We would like to use the K-means method to make the clustering automatically
 3077 and justify our above intuitive conclusion from the scatter diagram. This example,
 3078 although simple, is closer to reality. You can apply the procedures described here
 3079 to your own data for various clustering purposes, such as defining wildfire zones,
 3080 ecoregions, climate regimes, and agricultural areas. You can find numerous applica-
 3081 tion examples by an online search using keywords like “K-means clustering for
 3082 ecoregions”, or “K-means clustering for climate regimes.”

3083 The daily weather data for the Miami International Airport can be obtained from
 3084 the Global Historical Climatology Network-Daily (GHCN-D). The station code is
 3085 USW00012839. The daily data from 2001 to 2020 are included in the book’s master
 3086 data file `data.zip` from the book website www.climatestatistics.org. The file
 3087 name of the Miami data is

3088 `MiamiIntlAirport2001_2020.csv`

3089 You can access the updated data by an online search for “NOAA Climate Data
 3090 Online USW00012839.” Figure 3.2(b) shows the K-means clustering result for the
 3091 daily T_{min} vs WDF2 for the 2001 daily data at the Miami International Airport.
 3092 The two clusters generated by the K-means method agree with our intuition by
 3093 looking at the scatter diagram Fig. 3.2(a). However, the K-means clustering result
 3094 Fig. 3.2(b) provides a more clear picture: A red cluster with each data point indi-
 3095 cated by a small triangle, and a black cluster with each data point indicated by a
 3096 small circle. The cluster boundaries are linked by dashed lines. The large vertical
 3097 gap at T_{min} around 5°C (in winter) between the two clusters indicates that the
 3098 fastest 2-minute wind in a day has a WDF2 angle value around 50° , a northeast
 3099 wind for the red cluster, and around 300° , a northwest wind for the black cluster.

3100 Figure 3.2 may be generated by the following computer codes.

```
3101 #R for Fig. 9.2: K-means clustering for 2001 daily weather
3102 #data at Miami International Airport, Station ID USW00012839
3103 setwd("~/climstats")
3104 dat = read.csv("data/MiamiIntlAirport2001_2020.csv",
3105   header=TRUE)
3106 dim(dat)
3107 #[1] 7305 29
3108 tmin = dat[, 'TMIN']
3109 wdf2 = dat[, 'WDF2']
3110 # plot the scatter diagram Tmin vs WDF2
3111 setEPS() #Plot the data of 150 observations
3112 postscript("fig0902a.eps", width=5, height=5)
3113 par(mar=c(4.5, 4.5, 2, 4.5))
3114 plot(tmin[2:366], wdf2[2:366],
3115   pch =16, cex = 0.5,
3116   xlab = 'Tmin [degC]',
3117   ylab = 'Wind Direction [deg]', grid())
3118 title('a) 2001 Daily Miami Tmin vs WDF2',
3119   cex.main = 0.9, line = 1)
3120 axis(4, at = c(0, 45, 90, 135, 180, 225, 270, 315, 360),
3121   lab = c('N', 'NE', 'E', 'SE', 'S', 'SW', 'W', 'NW', 'N'))
3122 mtext('Wind Direction', side = 4, line =3)
```

```
3123 dev.off()
3124 #K-means clustering
3125 K = 2 #assuming K = 2, i.e., 2 clusters
3126 mydata = cbind(tmin[2:366], wdf2[2:366])
3127 fit = kmeans(mydata, K) # K-means clustering
3128 #Output the coordinates of the cluster centers
3129 fit$centers
3130 #1 18.38608 278.8608
3131 #2 21.93357 103.9161
3132 fit$tot.withinss # total WCSS
3133 #[1] 457844.9 for # the value may vary for each run
3134
3135 #Visualize the clusters by kmeans.ani()
3136 mycluster <- data.frame(mydata, fit$cluster)
3137 names(mycluster)<-c('Tmin[degC]',
3138                      'WindDirection[deg]',
3139                      'Cluster')
3140 library(animation)
3141 par(mar = c(4.5, 4.5, 2, 4.5))
3142 kmeans.ani(mycluster, centers = K, pch=1:K, col=1:K,
3143             hints = '')
3144 title(main=
3145         "(b) K-means Clusters for Daily Tmin vs WDF2",
3146         cex.main = 0.8)
3147 axis(4, at = c(0, 45, 90, 135, 180, 225, 270, 315, 360),
3148       lab = c('N', 'NE', 'E', 'SE', 'S', 'SW', 'W', 'NW', 'N'))
3149 mtext('WindDirection', side = 4, line =3)
```

```

# Python plot Fig. 9.2: K-means clustering for Tmin and WDF2
# K-means clustering for Fig. 9.2(b)
K = 2
# combine data into two column df
mydata = pd.concat([tmin[1:366],wdf2[1:366]], axis = 1)

# K-means clustering
fit = KMeans(n_clusters=K).fit(mydata)

# cluster center coordinates
print(fit.cluster_centers_, '\n')#enters of the two clusters
print(fit.inertia_) # total WCSS

#kmeans = KMeans(n_clusters=K, random_state=0)
mydata['cluster'] = \
kmeans.fit_predict(mydata[['TMIN', 'WDF2']])
colors = ['black', 'red']
mydata['color'] = \
mydata.cluster.map({0:colors[0], 1:colors[1]})

kmeans = KMeans(n_clusters=K, random_state=0)
mydata['cluster'] = \
kmeans.fit_predict(mydata[['TMIN', 'WDF2']])

colors = ['black', 'red']
mydata['color'] = \
mydata.cluster.map({0:colors[0], 1:colors[1]})

x1 = fit.cluster_centers_[0][0] #Cluster center
x2 = fit.cluster_centers_[1][0]
y1 = fit.cluster_centers_[0][1]
y2 = fit.cluster_centers_[1][1]

plt.title("(b) K-means for Tmin vs WDF2", pad = 10)
plt.xlabel(r" Tmin [$\degree C]", labelpad = 10)
plt.ylabel(r" Wind Direction [$\degree ]", labelpad = 10)
plt.yticks([0,50,100,200,300])
# plot points
plt.scatter(mydata['TMIN'], mydata['WDF2'],
            s = 10, color = mydata['color']) #s is size

for i in mydata.cluster.unique():
    points = \
        mydata[mydata.cluster == i][['TMIN', 'WDF2']].values
    # get convex hull
    hull = ConvexHull(points)
    # get x and y coordinates
    # repeat last point to close the polygon
    x_hull = np.append(points[hull.vertices,0],
                        points[hull.vertices,0][0])
    y_hull = np.append(points[hull.vertices,1],
                        points[hull.vertices,1][0])
    # plot shape
    plt.fill(x_hull, y_hull, alpha=0.2, c=colors[i])
    # plot centers
    plt.scatter([x1,x2],[y1,y2], color = ['red', 'black'],
                linewidth = 2, facecolors= 'none', s = 700)
plt.show()

```

3.1.4.2 Why $K = 2$ for the Miami weather data?

The K-means computer code for the Miami Tmin and WDF2 data shows two cluster centers at $C_1(18.4^\circ\text{C}, 278.9^\circ)$, and $C_2(21.9^\circ\text{C}, 103.9^\circ)$. The corresponding total WCSS is 457844.9. To justify our choice of $K = 2$, we compute tWCSS, the total WCSS, for different numbers of K , and then use the elbow rule to determine the best K value. Figure 3.3(a) shows the tWCSS scores for $K = 2, 3, \dots, 8$. We can see that the elbow appears at $K = 2$, since total WCSS decrement suddenly slows down from $K = 2$. Thus, the elbow rule has justified our choice of $K = 2$.

Another method to choose K is the knee rule, which is determined by the variances explained by K clusters. When $K = 1$, tWCSS is the spatial variance of the entire data. When $K = N$, tWCSS = 0, where N is the total number of the data points. Usually, we have

$$\text{tWCSS}[1] \geq \text{tWCSS}[K] \geq \text{tWCSS}[N] \quad (3.23)$$

for any $1 \leq K \leq N$. Then,

$$\text{pWCSS}[K] = \frac{\text{tWCSS}[1] - \text{tWCSS}[K]}{\text{tWCSS}[1]} \times 100\% \quad (3.24)$$

is defined as the percentage of variance explained by the K clusters. Usually, pWCSS[K] is an increasing function of K . When the increase rate suddenly slows down from K , a shape like a knee appears in the pWCSS[K] curve. We then use this K as the optimal number of centers for our K-means clustering. This is the so called the knee rule. Figure 3.3(b) shows that the increase of pWCSS[K] dramatically slows down at $K = 2$. Thus, we choose $K = 2$, which provides another justification why $K = 2$.

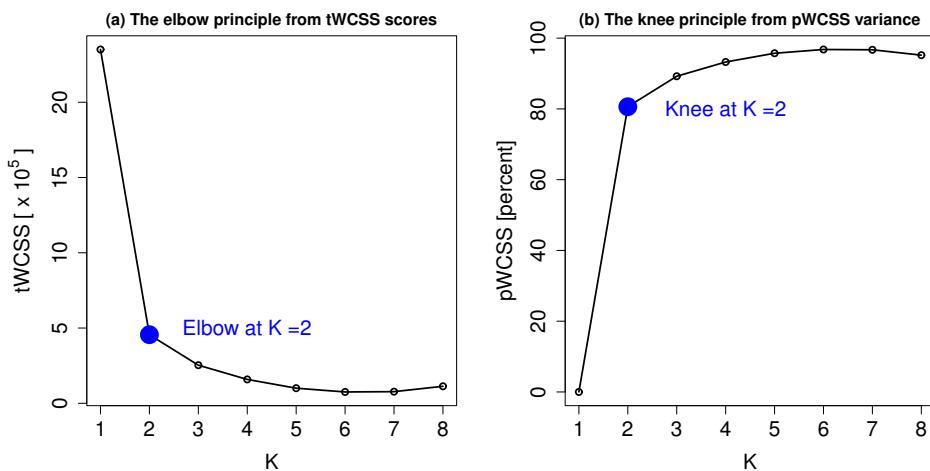


Fig. 3.3 (a) tWCSS scores for different K . (b) pWCSS variances for different K .

Figure 3.3 may be generated by the following computer codes.

```
3172 #R plot Fig. 9.3: tWCSS(K) and pWCSS(K)
3173 twcss = c()
3174 for(K in 1:8){
3175   mydata=cbind(tmax[2:366], wdf2[2:366])
3176   twcss[K] = kmeans(mydata, K)$tot.withinss
3177 }
3178 twcss
3179 par(mar = c(4.5, 6, 2, 0.5))
3180 par(mfrow=c(1,2))
3181 plot(twcss/100000, type = 'o', lwd = 2,
3182       xlab = 'K', ylab = bquote('tWCSS'[x] ~ 10^5 ~ ]),
3183       main = '(a) The elbow principle from tWCSS scores',
3184       cex.lab = 1.5, cex.axis = 1.5)
3185 points(2, twcss[2]/100000, pch =16, cex = 3, col = 'blue')
3186 text(4, 5, 'Elbow at K=2', cex = 1.5, col = 'blue')
3187 #compute percentage of variance explained
3188 pWCSS = 100*(twcss[1]- twcss)/twcss[1]
3189 plot(pWCSS, type = 'o', lwd = 2,
3190       xlab = 'K', ylab = 'pWCSS[percent]',
3191       main = '(b) The knee principle from pWCSS variance',
3192       cex.lab = 1.5, cex.axis = 1.5)
3193 points(2, pWCSS[2], pch =16, cex = 3, col = 'blue')
3194 text(4, 80, 'Knee at K=2', cex = 1.5, col = 'blue')
3195 dev.off()
```

```

#Python plot Fig. 9.3: tWCSS(K) and pWCSS(K)
#Plot Fig. 9.3(a): Elbow principle
twcss = np.zeros(9)
for K in range(1,9):
    mydata = pd.concat([tmin[1:366],wdf2[1:366]], axis = 1)
    twcss[K] = KMeans(n_clusters=K).fit(mydata).inertia_
# remove K = 0 value from twcss
twcss = twcss[1:]
# plot elbow
plt.plot(np.linspace(1,8,8), twcss/100000, 'k', linewidth=2)
# plot points
plt.scatter([np.linspace(1,8,8)], [twcss/100000], color='k')

# plot elbow at K = 2
plt.scatter(2,twcss[1]/100000, color = 'blue', s = 500)
# add text
plt.text(2.2, 5.2, 'Elbow at K=2', color='blue', size=20)
# add labels
plt.title("(a) The elbow principle from tWCSS scores",
          pad = 10)
plt.xlabel("K", labelpad = 10)
plt.ylabel(r"tWCSS [x $10^5$]", labelpad = 10)
plt.show()

#Plot Fig. 9.3(b): Knee principle
# compute percentage of variance explained
pWCSS = 100*(twcss[0]-twcss)/twcss[0]
# plot knee
plt.plot(np.linspace(1,8,8),pWCSS, 'k', linewidth = 2)
# add points
plt.scatter(np.linspace(1,8,8),pWCSS, color = 'k')

# plot knee at K = 2
plt.scatter(2,pWCSS[1], color = 'blue', s = 500)
# add text
plt.text(2.2, 76, 'Knee at K=2', color = 'blue',
         size = 20)
# add labels
plt.title("(b) The knee principle from pWCSS variance",
          pad = 10)
plt.xlabel("K", labelpad = 10)
plt.ylabel("pWCSS [%]", labelpad = 10)
plt.show()

```

3196

3.1.4.3 Steps of the K-means clustering data analysis

3198 From the previous sub-sections, we may summarize that the K-means method for
 3199 data analysis consists of the following steps.

- 3200 (i) Data preparation: Write your data into an $N \times 2$ matrix, where N is the total
 3201 number of data points, and 2 is for two variables, i.e., in a 2-dimensional
 3202 space. The real value matrix must have no missing data. Although our

3203 previous examples are for 2-dimensional data, the K-means clustering can
 3204 be used for higher dimensional data. On the Internet, you can easily find
 3205 examples of the K-means for data of three or more variables.

- 3206 (ii) Preliminary data analysis: You may wish to make some preliminary simple
 3207 analyses before applying the K-means method. These analyses may include
 3208 plotting a scatter diagram, plot a histogram for each variable, and compute
 3209 mean, standard deviation and quantiles for each variable. If the data for
 3210 the two variables have different units, you may convert the data into their
 3211 standardized anomalies. The standardized anomaly computing procedure
 3212 is called “scale” in the machine learning community. For example, the R
 3213 command `scale(x)` yields the standard anomalies of data sequence `x`, i.e,
 3214 divided by the standard deviation after a subtraction of mean. Mathemati-
 3215 cally speaking, the scaling procedure is preferred when two variables have
 3216 different units, because the $\text{unit1}^2 + \text{unit2}^2$ in tWCSS usually do not have a
 3217 good interpretation. In contrast, standardized anomalies are dimensionless
 3218 and have no units. Non-dimensional data are often preferred when applying
 3219 statistical methods. However, our goal of the K-means clustering is to find
 3220 patterns without labeling the data, which is an unsupervised learning. The
 3221 K-means results from scaled data and unscaled data may not be the same.
 3222 As long as the K-means result is reasonable and interpretable, regardless
 3223 of the scaled or unscaled data, we adopt the result. For example, in our
 3224 Miami Tmin and WDF2 example in the previous sub-section, we found the
 3225 K-means clusters for both scaled and unscaled data, but we have chosen
 3226 the result from the unscaled data from our view point of interpretation.
 3227 Therefore, we keep in mind that the K-means is a result-oriented method.
- 3228 (iii) Test K-means clustering: Apply a computer code of K-means to the prepared
 3229 $N \times 2$ data matrix with a K value inspired from the scatter diagram or
 3230 science background of the data. You may apply the code once to the scaled
 3231 data and another to the unscaled data, and see if the results make sense.
 3232 Run your code a few times and see if the results vary.
- 3233 (iv) Determine the optimal K by the elbow rule: Compute the tWCSS scores for
 3234 different K and find the elbow location of the tWCSS(K) curve. Apply the
 3235 K-means computer code for this optimal K value. Compare this K with
 3236 the K in the previous step. Check the details of the clustering results, such
 3237 as cluster assignment, tWCSS, centers.
- 3238 (v) Visualize the K-means clustering results: Use a computer code to plot the K
 3239 clusters. Different visualization codes are available for your choice.
- 3240 (vi) Interpret the final result: Identify the characteristics of the clusters and inter-
 3241 pret them with text or more figures.

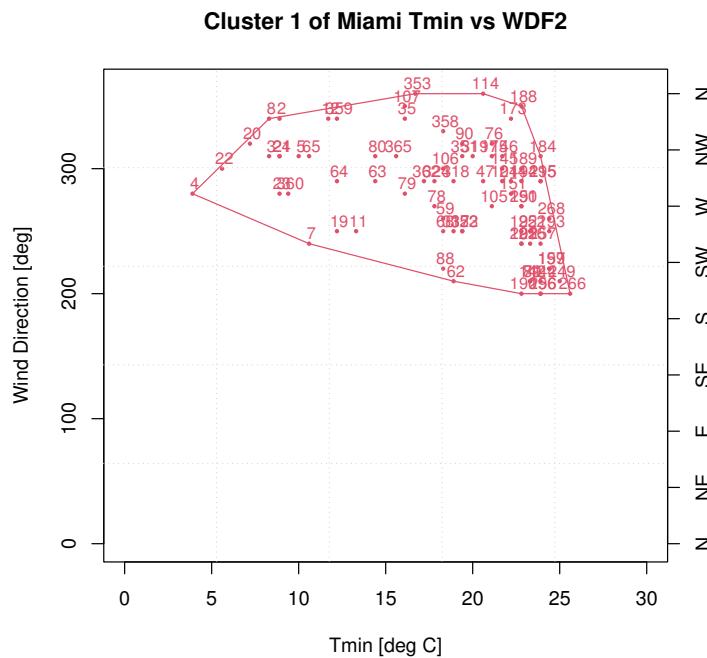
3242 Applying the K-means computer code is very easy, but in the K-means data
 3243 analysis, we should avoid the following common mistakes:

- 3244 • Skipped the preliminary data analysis step: Some people may quickly apply the

3245 K-means computer code and treat the K-means results as unique, hence regard
 3246 their their plot of the clusters as the end of the analysis.

- 3247 • Missed the interpretation: Interpretation step is needed, short or long. This is
 3248 our purpose anyway.
- 3249 • Forgot to scale the data: We should always analyze the scaled data. At the same
 3250 time, because of the interpretation requirement, we should also check the re-
 3251 sults from the unscaled data, and see which result is more reasonable.
- 3252 • Overlooked the justification of the optimal K : A complete K-means clustering
 3253 should have this step, which helps us to explore all the possible clusters.

3254 In addition to using the standard K-means computer code you can find from the
 3255 Internet or that built in R or Python, you sometimes may wish to make special
 3256 plots for your clustering results with some specific features, such as a single cluster
 3257 on a figure with the marked data order, as shown in Fig. 3.4. Convex hull is a
 3258 simple method to plot a cluster. By definition, the convex hull of a set of points is
 3259 the smallest convex polygon that encloses all of the points in the set. Below is the
 3260 a computer code to plot Fig. 3.4 using the convex hull method.

**Fig. 3.4**

Convex hull for Cluster 1. The number above each data point is the day number in a year, e.g., January 1 being 1, February 2 being 33, and December 31 being 365 in 2001. The day numbers help identify the seasonality of the points, such as, the left most tip of the cluster being Day 4, a specific winter day.

3261 Figure 3.4 may be generated by the following computer codes.

```

3262 #R plot Fig. 9.4: Convex hull for a cluster
3263 setwd("~/climstats")
3264 dat = read.csv("data/MiamiIntlAirport2001_2020.csv",
3265                 header=TRUE)
3266 dim(dat)
3267 #[1] 7305   29
3268 tmin = dat[, 'TMIN']
3269 wdf2 = dat[, 'WDF2']
3270 #K-means clustering
3271 K = 2 #assuming K = 2, i.e., 2 clusters
3272 mydata = cbind(tmin[2:366], wdf2[2:366]) #2001 data
3273 clusterK = kmeans(mydata, K) # K-means clustering
3274 mycluster <- data.frame(mydata, clusterK$cluster)
3275 plotdat = cbind(1:N, mycluster)
3276
3277 par(mar = c(4.5, 6, 2, 0.5)) #set up the plot margin
3278 i = 1 # plot Cluster 1
3279 N = 365 #Number of data points
3280 X = plotdat[which(mycluster[,3] == i), 1:3]
3281 colnames(X)<-c('Day', 'Tmin[degC]', 'WDF2[deg]')
3282 plot(X[,2:3], pch = 16, cex = 0.5, col = i,
3283       xlim = c(0, 30), ylim = c(0, 365),
3284       main = 'Cluster 1 of Miami Tmin vs WDF2')
3285 grid(5, 5)
3286 #chull() finds the boundary points of a convex hull
3287 hpts = chull(X[, 2:3])
3288 hpts1 = c(hpts, hpts[1]) #close the boundary
3289 lines(X[hpts1, 2:3], col = i)
3290 for(j in 1:length(X[,1])){
3291   text(X[j,2], X[j,3] + 8, paste("", X[j,1]),
3292         col = i, cex = 0.8)
3293 } #Put the data order on the cluster

```

```

# Python plot Fig. 9.4: Convex hull for a cluster
K = 2
clusterK = KMeans(n_clusters=K).fit(mydata)
mydata['cluster'] = \
    kmeans.fit_predict(mydata[['TMIN', 'WDF2']])
mydata['day'] = np.arange(1, 366)
subset = mydata[mydata['cluster']==0]
# set up plot
plt.title("Cluster 1 of Miami Tmin vs WDF2", pad = 10)
plt.xlabel(r"Tmin [$\degree C]", labelpad = 10)
plt.ylabel(r"WDF2 [$\degree ]", labelpad = 10)

plt.xlim([0,30])
plt.ylim([0,400])
plt.yticks([0,100,200,300])

# plot points
plt.scatter(subset["TMIN"], subset["WDF2"], color = 'k')
# add labels
for i in range(len(subset["TMIN"])):
    plt.text(subset["TMIN"].values[i],
             subset["WDF2"].values[i],
             subset["day"].values[i], size=12)
# add boundary
for i in subset.cluster.unique():
    points = \
        subset[subset.cluster == i][['TMIN', 'WDF2']].values
    # get convex hull
    hull = ConvexHull(points)
    # get x and y coordinates
    # repeat last point to close the polygon
    x_hull = np.append(points[hull.vertices,0],
                        points[hull.vertices,0][0])
    y_hull = np.append(points[hull.vertices,1],
                        points[hull.vertices,1][0])
    # plots shape
    plt.fill(x_hull, y_hull, alpha=0.2, c=colors[i])

```

3294

3295 The K-means method can not only be used for weather data clustering and
 3296 climate classification, but also for many other purposes, such as weather forecasting,
 3297 storm identification, and more. In many applications, K-means is part of a machine
 3298 learning algorithm, and takes a modified version, such as the incremental K-means
 3299 for weather forecasting and the weighted K-means for climate analysis.

3300

3.2 Support vector machine

3301

3302 The K-means clustering can organize a suite of given data points into K clusters.
 3303 The data points are not labeled. This is like organizing basket of fruits into K piles

according to the minimum tWCSS principle. After the clustering, we may name the clusters, say $1, 2, \dots, K$, or apples, oranges, peaches, and grapes. ML often uses the term “cluster label” in lieu of “cluster name.” So, we say that we label a cluster when we name a cluster. The learning process for a un-labeled data is called unsupervised learning. If a basket of fruits is sorted out by a baby who cannot tell the fruit names according to wishes of his mother, he is performing an unsupervised learning.

In contrast, support vector machine (SVM) is in general a supervised learning, working on the labelled data, e.g., a basket of apples, oranges, peaches and grapes. Based on the numerical data, SVM has two purposes. First, SVM determines the maximum “distances” among the sets of labeled data, say apples, oranges, peaches and grapes. This is a training process and helps a computer learn the differences among different labeled sets, e.g., different fruits, different weather, or different climate regimes. Second, SVM makes a prediction, which, based on the trained model, predicts the labels (i.e., names) of the new and unlabeled data. Thus, SVM is a system to learn the maximum differences among the labelled clusters of data, and to predict the labels of the new unlabeled data. This process of training and prediction mimics a human’s learning experience. For example, through years a child is trained by his parents to learn the knowledge of apples, oranges and other fruits. At a certain time, the child is well trained and ready to independently tell (i.e., predict) whether a new fruit is apple or orange. During this learning process, the fruits are labeled (e.g., apples, oranges, etc). His parents help input the data into his brain through a teaching and learning process. His brain was trained to maximize the difference between an apple and an orange. In the prediction stage, the child independently imports the data using his eyes, and make the prediction whether a fruit is an apple or orange. SVM mimics this process by maximizing the difference between two or more labeled categories or clusters in the training process, then by using the trained model, SVM predicts the categories for the new data. The key is how to quantify the difference between any two categories.

In the following we will illustrate these SVM learning ideas using three examples:
 (i) A trivial case of two data points, (ii) training and prediction for a system of three data points, and (iii) a more general case of many data points.

3.2.0.1 Maximize the difference between two labelled points

We formulate a two-point learning question as follows. We have an apple and an orange. We know all the data about the apple and orange. We teach a baby how to use the data to quantify the maximum difference between the apple and orange. In the machine learning language, “apple” and “orange” are called labels. The data are called the training data, such as color and surface smoothness. The maximum difference will be used to predict whether a new fruit is an apple or orange.

An abstract description of the above may be put in the following way. Given the data for an object labeled A and also the data for another object labeled B, quantify the maximum difference between A and B. A simple example corresponding to this

statement is that given two points $P_1(1, 1)$ and $P_2(3, 3)$ which are labeled $y_1 = -1$ and $y_2 = 1$, quantify the maximum difference between y_1 and y_2 . For this simple problem, the obvious answer is that the two categories labeled by y_1 and y_2 are divided by the perpendicular bisector of the two points P_1 and P_2 as shown in Fig. 3.5. The bisector is called the *separating hyperplane*. A hyperplane means a plane in a space of more than three dimensions. In a 2-dimensional space, it is a straight line. In a 3-dimensional space, it is just a regular flat plane. The real ML applications are often for higher dimensional data, e.g., a comprehensive description of an apple may need data of color, surface smoothness, stem size, skin thickness, and more, which form a p -dimensional vector.

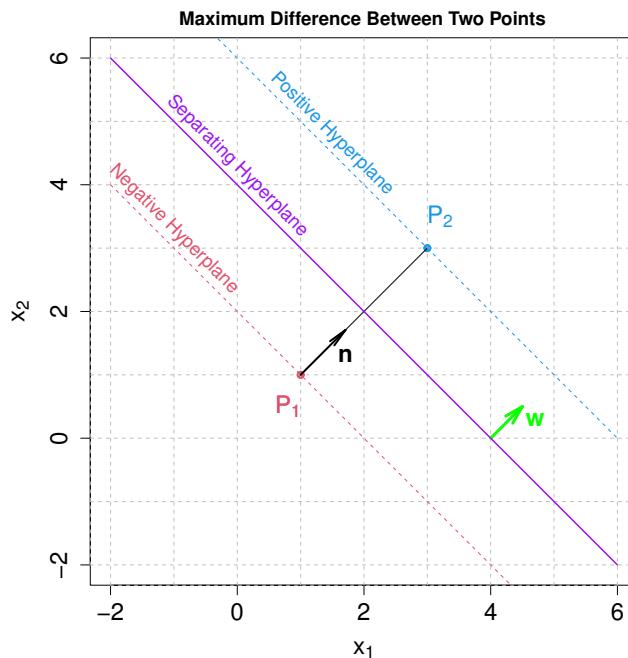


Fig. 3.5 Quantify the maximum difference between two points.

The maximum margin of difference between the two categories are bounded by two lines parallel to the bisector, one line passing through P_1 and another P_2 . One line is called the negative hyperplane, and another the positive hyperplane. The data points on these hyperplanes are called *support vectors*. The support vector machine (SVM) has its name from these vectors.

The equation for the separating hyperplane (i.e., the purple line) is

$$x_1 + x_2 = 4, \quad (3.25)$$

that for the positive hyperplane (i.e., the blue dashed line) is

$$x_1 + x_2 = 6, \quad (3.26)$$

³³⁶³ and that for the negative hyperplane (i.e., the red dashed line) is

$$x_1 + x_2 = -2. \quad (3.27)$$

³³⁶⁴ The unit vector in the perpendicular bisector direction from P_1 to P_2 is

$$\mathbf{n} = \left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right). \quad (3.28)$$

³³⁶⁵ This vector is a normal vector for the separating hyperplane. The equation of the
³³⁶⁶ separating hyperplane may be expressed in the normal vector form

$$\mathbf{n} \cdot \mathbf{x} = 2\sqrt{2}. \quad (3.29)$$

³³⁶⁷ This can further be written in the following form

$$\mathbf{w} \cdot \mathbf{x} - b = 0, \quad (3.30)$$

³³⁶⁸ where

$$\mathbf{w} = (0.5, 0.5) \quad (3.31)$$

³³⁶⁹ and

$$b = 2. \quad (3.32)$$

³³⁷⁰ An advantage of this form is that the equations for the positive and negative hy-
³³⁷¹ perplanes can be simply written in the form

$$\mathbf{w} \cdot \mathbf{x} - b = \pm 1, \quad (3.33)$$

³³⁷² and that for the separating hyperplane can be written as

$$\mathbf{w} \cdot \mathbf{x} - b = 0. \quad (3.34)$$

³³⁷³ In the above, \mathbf{w} is called a normal vector and is related with the unit normal vector
³³⁷⁴ \mathbf{n} as follows:

$$\mathbf{w} = -\frac{1}{\sqrt{2}} \mathbf{n}. \quad (3.35)$$

³³⁷⁵ The distance between the positive and negative hyperplanes is

$$D_m = \frac{2}{|\mathbf{w}|} = 2\sqrt{2}, \quad (3.36)$$

³³⁷⁶ is the margin size. The maximum difference between P_1 and P_2 is quantified by
³³⁷⁷ D_m . The primary goal of SVM is to maximize D_m when training data are given.

³³⁷⁸ The distance of the origin to the separating hyperplane is

$$D_o = -\frac{b}{|\mathbf{w}|} = 2\sqrt{2}. \quad (3.37)$$

³³⁷⁹ These formulations of the hyperplanes and the quantities \mathbf{w}, b, D_m can be used
³³⁸⁰ for the general SVM mathematical formulation, which are described in the next
³³⁸¹ sub-section.

³³⁸² Figure 3.5 may be generated by the following computer code.

```

3383 #R plot Fig. 9.5: Maximum difference between two points
3384 x = matrix(c(1, 1, 3, 3),
3385   ncol = 2, byrow = TRUE)#Two points
3386 y= c(-1, 1) #Two labels -1 and 1
3387 #Plot the figure and save it as a .eps file
3388 setEPS()
3389 postscript("fig0905.eps", height=7, width=7)
3390 par(mar = c(4.5, 4.5, 2.0, 2.0))
3391 plot(x, col = y + 3, pch = 19,
3392   xlim = c(-2, 6), ylim = c(-2, 6),
3393   xlab = bquote(x[1]), ylab = bquote(x[2]),
3394   cex.lab = 1.5, cex.axis = 1.5,
3395   main = "Maximum\u20d7Difference\u20d7Between\u20d7Two\u20d7Points")
3396 axis(2, at = (-2):6, tck = 1, lty = 2,
3397   col = "grey", labels = NA)
3398 axis(1, at = (-2):6, tck = 1, lty = 2,
3399   col = "grey", labels = NA)
3400 segments(1, 1, 3, 3)
3401 arrows(1, 1, 1.71, 1.71, lwd = 2,
3402   angle = 9, length= 0.2)
3403 text(1.71, 1.71-0.4, quote(bold('n')), cex = 1.5)
3404 arrows(4, 0, 4 + 0.5, 0 + 0.5, lwd = 3,
3405   angle = 15, length= 0.2, col = 'green' )
3406 text(4.7, 0.5 -0.2, quote(bold('w')), cex = 1.5, col = 'green')
3407 x1 = seq(-2, 6, len = 31)
3408 x20 = 4 - x1
3409 lines(x1, x20, lwd = 1.5, col = 'purple')
3410 x2m = 2 - x1
3411 lines(x1, x2m, lty = 2, col = 2)
3412 x2p = 6 - x1
3413 lines(x1, x2p, lty = 2, col = 4)
3414 text(1-0.2,1-0.5, bquote(P[1]), cex = 1.5, col = 2)
3415 text(3+0.2,3+0.5, bquote(P[2]), cex = 1.5, col = 4)
3416 text(1-0.2,1-0.5, bquote(P[1]), cex = 1.5, col = 2)
3417 text(0,4.3, 'Separating\u20d7Hyperplane',
3418   srt = -45, cex = 1.2, col = 'purple')
3419 text(1.5, 4.8, 'Positive\u20d7Hyperplane',
3420   srt = -45, cex = 1.2, col = 4)
3421 text(-1, 3.3, 'Negative\u20d7Hyperplane',
3422   srt = -45, cex = 1.2, col = 2)
3423 dev.off()

```

```

#Python plot Fig. 9.5: Maximum difference between two points
# define the two points
x = np.array([1,1,3,3]).reshape(2,2)
# define the two labels
y = (-1,1)
# set up plot
plt.figure(figsize = (10,10))
plt.title("Maximum Difference Between Two Points", pad = 10)
plt.xlabel("$x_1$", labelpad = 10)
plt.ylabel("$x_2$", labelpad = 10)
plt.xlim(-2.2,6.2)
plt.xticks([-2,-1,0,1,2,3,4,5,6],[-2,'',0,'',2,'',4,'',6])
plt.ylim(-2.2,6.2)
plt.yticks([-2,-1,0,1,2,3,4,5,6],[-2,'',0,'',2,'',4,'',6])
plt.grid()
# plot points
plt.scatter(x[0],x[1], color = 'r')
plt.scatter(x[2],x[3], color = 'dodgerblue')
# connect points
plt.plot([1,3],[1,3], 'k')
# add dashed lines
plt.plot([0,6],[6,0], color = 'dodgerblue', linestyle = '--')
plt.plot([-2,6],[6,-2], color = 'purple', linestyle = '--')
plt.plot([-2,4],[4,-2], 'r', linestyle = '--')
# plot arrows
plt.arrow(1,1,0.5,0.5, head_width = 0.12, color='k')
plt.arrow(4,0,0.3,0.3, head_width = 0.12, color='limegreen')
# add text
plt.text(0.5,4.3, "Positive Hyperplane", rotation = -45,
         size = 14, color = 'dodgerblue')
plt.text(3.2,3.2, "$P_2$", size = 15, color = 'dodgerblue')
plt.text(-1.2,3.8, "Separating Hyperplane", rotation = -45,
         size = 14, color = 'purple')
plt.text(4.4,0.1, "w", size = 15, color = 'limegreen')
plt.text(1.5,1.2, "n", size = 15)
plt.text(-2,2.7, "Negative Hyperplane", rotation = -45,
         size = 14, color = 'r')
plt.text(0.7, 0.6, "$P_1$", size = 15, color = 'r')
plt.show()

```

3425

3.2.1 SVM for a system of three points labeled in two categories

Let us progress from a two-point system to a three-point system and build an SVM using a computer code. Figure 3.6 shows the three points labeled in two categories. The two blue points are $P_1(1,1), P_2(2,1)$ which are in the first category corresponding to the categorical value $y = 1$. The red point $P_3(3,3.5)$ is in the second category corresponding to $y = -1$. Training these data is to maximize the margin size D_m , which results in the following SVM parameters.

$$w = (-0.28 - 0.69), \quad b = 2.24, \quad D_m = 2.69. \quad (3.38)$$

3434 The support vectors are $P_2(2, 1)$ and $P_3(3, 3.5)$.

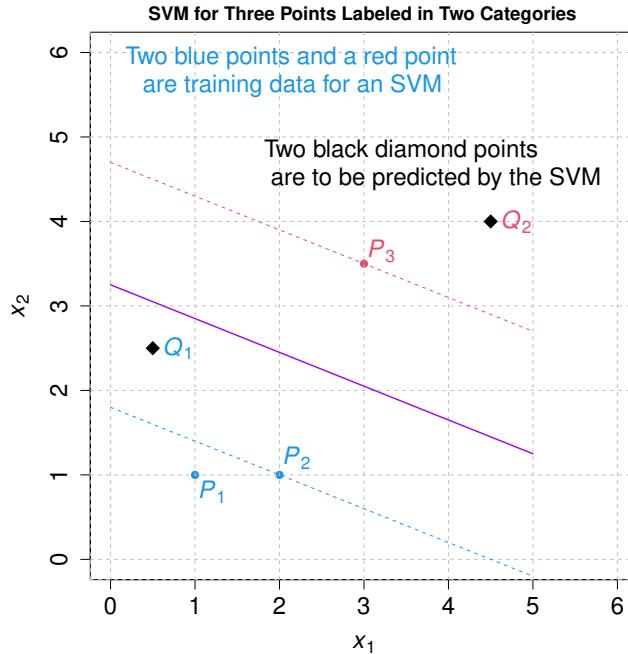


Fig. 3.6 SVM training from three data points, and the prediction by the trained SVM.

3435 We acquire two new data points $Q_1(0.5, 2.5)$ and $Q_2(4.5, 4.0)$. We wish to use
3436 the trained SVM to predict whether $Q_1(0.5, 2.5)$ corresponds to 1 (i.e., the first
3437 category) or -1 (i.e., the second category). The SVM prediction results are shown
3438 in Fig. 3.6: $Q_1(0.5, 2.5)$ belongs to the first category, and $Q_2(4.5, 4.0)$ to the second
3439 category. Geometrically, this prediction result in Fig. 3.6 is obvious and has nothing
3440 surprising. The significance of this example is its systematic use of the SVM method
3441 and its computer algorithm. Although this system is simple and has only three
3442 points, solving this problem by hand is very difficult. The following computer code
3443 solves this problem and generates Fig. 3.6.

```

3444 #R plot Fig. 9.6: SVM for three points
3445 #training data x
3446 x = matrix(c(1, 1, 2, 1, 3, 3.5),
3447   ncol = 2, byrow = TRUE)
3448 y = c(1, 1, -1) #two categories 1 and -1
3449 plot(x, col = y + 3, pch = 19,
3450   xlim = c(-2, 8), ylim = c(-2, 8))
3451 library(e1071)
3452 dat = data.frame(x, y = as.factor(y))
3453 svm3P = svm(y ~ ., data = dat,
3454   kernel = "linear", cost = 10,
3455   scale = FALSE,
3456   type = 'C-classification')
3457 svm3P #This is the trained SVM

```

```

3458 xnew = matrix(c(0.5, 2.5, 4.5, 4),
3459   ncol = 2, byrow = TRUE) #New data
3460 predict(svm3P, xnew) #prediction using the trained SVM
3461 # 1 2
3462 # 1 -1
3463
3464 # Find hyperplane, normal vector, and SV (wx + b = 0)
3465 w = t(svm3P$coefs) %*% svm3P$SV
3466 w
3467 #[1,] -0.2758621 -0.6896552
3468 b = svm3P$rho
3469 b
3470 #[1] -2.241379
3471 2/norm(w, type ='2') #maximum margin of separation
3472 #[1] 2.692582
3473
3474 x1 = seq(0, 5, len = 31)
3475 x2 = (b - w[1]*x1)/w[2]
3476 x2p = (1 + b - w[1]*x1)/w[2]
3477 x2m = (-1 + b - w[1]*x1)/w[2]
3478 x20 = (b - w[1]*x1)/w[2]
3479 #plot the SVM results
3480 setEPS()
3481 postscript("fig0906.eps", height=7, width=7)
3482 par(mar = c(4.5, 4.5, 2.0, 2.0))
3483 plot(x, col = y + 3, pch = 19,
3484   xlim = c(0, 6), ylim = c(0, 6),
3485   xlab = bquote(x[1]), ylab = bquote(x[2]),
3486   cex.lab = 1.5, cex.axis = 1.5,
3487   main = 'SVM for three points labeled in two categories')
3488 axis(2, at = (-2):8, tck = 1, lty = 2,
3489   col = "grey", labels = NA)
3490 axis(1, at = (-2):8, tck = 1, lty = 2,
3491   col = "grey", labels = NA)
3492 lines(x1, x2p, lty = 2, col = 4)
3493 lines(x1, x2m, lty = 2, col = 2)
3494 lines(x1, x20, lwd = 1.5, col = 'purple')
3495 xnew = matrix(c(0.5, 2.5, 4.5, 4),
3496   ncol = 2, byrow = TRUE)
3497 points(xnew, pch = 18, cex = 2)
3498 for(i in 1:2){
3499   text(xnew[i,1] + 0.5, xnew[i,2] , paste('Q',i),
3500     cex = 1.5, col = 6-2*i)
3501 }
3502 text(2.2,5.8, "Two blue points and a red point
3503 are training data for an SVM",
3504   cex = 1.5, col = 4)
3505 text(3.5,4.7, "Two black diamond points
3506 are to be predicted by the SVM",
3507   cex = 1.5)
3508 dev.off()

```

```

#Python plot Fig. 9.6: SVM for three points
# define x for ease of plotting
x = np.array([1,1,2,1,3,3.5])
# redefine x for svm
X = [[1, 1],[2,1],[3,3.5]]
y = [1, 1, -1]
# define SVM
svm3P = svm.SVC(kernel='linear')
# train SVM
svm3P.fit(X, y)
# predict new data using the trained SVM
print('The prediction is:', 
      svm3P.predict([[0.5, 2.5],[4.5,4]]), '\n')
# find hyperplane, normal vector, and SV (wx + b = 0)
w = svm3P.coef_[0]
print('w is:', w, '\n')
b = -svm3P.intercept_ #intercept
print('b is:', b, '\n')
# find the maximum margin of separation
print(2/np.linalg.norm(w))
x1 = np.linspace(0,5,5)
x2 = (b - w[0]*x1)/w[1]
x2p = (1 + b - w[0]*x1)/w[1]
x2m = (-1 + b - w[0]*x1)/w[1]
x20 = (b - w[0]*x1)/w[1]
# set up plot
plt.figure(figsize = (10,10))
plt.title("SVM for three points labeled in two categories",
          pad = 10)
plt.xlabel("$x_1$", labelpad = 10)
plt.ylabel("$x_2$", labelpad = 10)
plt.xlim(-0.2,5.2)
plt.ylim(-0.2,6.2)
plt.grid()
# plot original points
plt.scatter([x[::2]], [x[1::2]],
            color=['dodgerblue','dodgerblue','red'])
# add lines
plt.plot(x1,x2p, color = 'dodgerblue', linestyle = '--')
plt.plot(x1, x2m, color = 'red', linestyle = '--')
plt.plot(x1, x20, color = 'purple')
# plot diamonds
plt.scatter([0.5,4.5], [2.5,4], color = 'k', marker = 'D')
# add text
plt.text(0.65,2.45,"$Q_1$", size = 16, color = 'dodgerblue')
plt.text(4.65,3.95, "$Q_2$", size = 16, color = 'r')
plt.text(0.3,5.4, "Two blue points and a red point\nare training data for an SVM",
         size = 16, color = 'dodgerblue')
plt.text(1.7,4.5, "Two black diamond points\nare to be predicted by the SVM",
         size = 16, color = 'k')
plt.show()

```

3510 **3.2.2 SVM mathematical formulation for a system of many
3511 points in two categories**
3512

3513 The previous quantification of maximum separation between two points or three
3514 points can be generalized to the separation of many points of two categories labeled
3515 $(1, -1)$, or $((1, 2)$, or (A, B) , or another way. What is a systematic formulation of
3516 mathematics to make the best separation, by maximizing the margin? How can we
3517 use a computer code to implement the separation? This sub-section attempts to
3518 answer these questions.

3519 We have n training data points:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), \quad (3.39)$$

3520 where x_i is a p -dimensional vector (i.e., using p parameters to describe every data
3521 point in a category), and y_i is equal to 1 or -1 , a category indicator, $i = 1, 2, \dots, n$.

3522 The separating hyperplane has its linear equation as follows

$$\mathbf{w} \cdot \mathbf{x} - b = 0. \quad (3.40)$$

3523 Our SVM algorithm is to find the p -dimensional normal vector \mathbf{w} and a real scalar
3524 b so that the distance

$$D_m = \frac{2}{|\mathbf{w}|} \quad (3.41)$$

3525 between the positive and negative hyperplanes

$$\mathbf{w} \cdot \mathbf{x} - b = \pm 1 \quad (3.42)$$

3526 is maximized.

3527 The maximization is the domain in the p -dimensional space below the negative
3528 hyperplane and above the positive hyperplane, i.e.,

$$\mathbf{w} \cdot \mathbf{x}_i - b \geq 1 \quad \text{when } y_i = 1 \quad (3.43)$$

3529 and

$$\mathbf{w} \cdot \mathbf{x}_i - b \leq -1 \quad \text{when } y_i = -1. \quad (3.44)$$

3530 This definition of the maximization implies that no training data points are located
3531 between the positive and negative hyperplanes. Thus, the solution of the maximization
3532 of $D_m = 2/|\mathbf{w}|$ must occur at the boundary of the domain, i.e., the positive
3533 and negative hyperplanes. The points on these hyperplanes are called the support
3534 vectors (SV). R or Python SVM code can find these SVs.

3535 Figure 3.7 shows an example of an SVM training with 20 training data points,
3536 and an SVM prediction for 3 data points. The data are given below:

| | | | | |
|------|----|----|-----|---|
| 3537 | X1 | X2 | y | |
| 3538 | 1 | 1 | 6.0 | 1 |
| 3539 | 2 | 2 | 8.0 | 1 |
| 3540 | 3 | 3 | 7.5 | 1 |

```

3541 4 1 8.0 1
3542 5 4 9.0 1
3543 6 5 9.0 1
3544 7 3 7.0 1
3545 8 5 9.0 1
3546 9 1 5.0 1
3547 10 5 3.0 2
3548 11 6 4.0 2
3549 12 7 4.0 2
3550 13 8 6.0 2
3551 14 9 5.0 2
3552 15 10 6.0 2
3553 16 5 0.0 2
3554 17 6 5.0 2
3555 18 8 2.0 2
3556 19 2 2.0 2
3557 20 1 1.0 2

```

3558 The new data to be predicted are

```

3559 [1,] 0.5 2.5
3560 [2,] 7.0 2.0
3561 [3,] 6.0 9.0

```

3562 The following computer code can do the SVM calculation and plot Fig. 3.7.

3563 We acquire three new data points $Q_1(0.5, 2.5)$, $Q_2(7.5, 2.0)$ and $Q_3(6.0, 9.0)$. We
 3564 wish to use the trained SVM to find out which point belongs to what category. The
 3565 SVM prediction result is shown in Fig. 3.7: $Q_1(0.5, 2.5)$ and $Q_2(7.5, 2.0)$ belong to
 3566 Category 1, and $Q_3(6.0, 9.0)$ to Category 2. Again, this prediction result is obvious
 3567 to us when the new data points are plotted, our eyes can see the points, and
 3568 our brains can make the decisions based on our trained intelligence. The SVM
 3569 applications in the real world often involve data of higher dimensions, in which
 3570 case the visual decision is impossible, and requires us to predict the category for
 3571 the new data without plotting. The prediction is based on the training data, new
 3572 data, and SVM algorithm without visualization. For example, the R prediction of
 3573 this problem can be done by the command `predict(svmP, xnew)` where `xnew` is
 3574 the new data in a 3×2 matrix, and `svmP` is the trained SVM. This command is
 3575 included in the computer code for generating Fig. 3.7.

```

3576 #R plot Fig. 9.7: SVM for many points
3577 #Training data x and y
3578 x = matrix(c(1, 6, 2, 8, 3, 7.5, 1, 8, 4, 9, 5, 9,
3579           3, 7, 5, 9, 1, 5,
3580           5, 3, 6, 4, 7, 4, 8, 6, 9, 5, 10, 6,
3581           5, 0, 6, 5, 8, 2, 2, 2, 1, 1),
3582           ncol = 2, byrow = TRUE)
3583 y= c(1, 1, 1, 1, 1, 1,

```

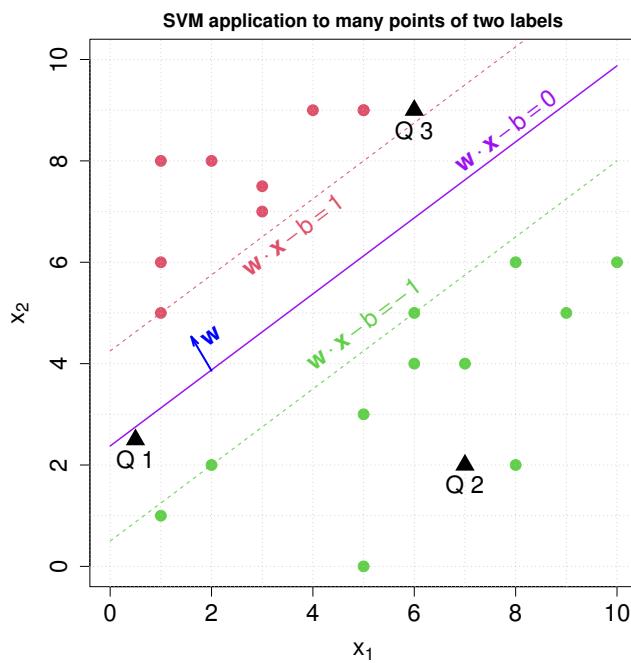


Fig. 3.7 SVM training from 20 data points, and the prediction of 3 new data points by the trained SVM.

```

3584      1,  1,  1,
3585      2,  2,  2,  2,  2,  2 ,
3586      2,  2,  2,  2)
3587
3588 library(e1071)
3589 dat = data.frame(x, y = as.factor(y))
3590 svmP = svm(y ~ ., data = dat,
3591             kernel = "linear", cost = 10,
3592             scale = FALSE,
3593             type = 'C-classification')
3594 svmP
3595 #Number of Support Vectors:  3
3596 svmP$SV #SVs are #x[9,], x[17,], x[19,]
3597 #9  1  5
3598 #17  6  5
3599 #19  2  2
3600
3601 # Find SVM parameters: w, b, SV (wx+c=0)
3602 w = t(svmP$coefs) %*% svmP$SV
3603 # In essence this finds the hyper plane that separates our points
3604 w
3605 #[1,] -0.39996  0.53328
3606 b = svmP$rho
3607 b
3608 #[1] 1.266573
3609 2/norm(w, type = '2')

```

```

3610 #[1] 3.0003 is the maximum margin
3611 x1 = seq(0, 10, len = 31)
3612 x2 = (b - w[1]*x1)/w[2]
3613 x2p = (1 + b - w[1]*x1)/w[2]
3614 x2m = (-1 + b - w[1]*x1)/w[2]
3615 x20 = (b - w[1]*x1)/w[2]
3616
3617 #plot the svm results
3618 setEPS()
3619 postscript("fig0907.eps", height=7, width=7)
3620 par(mar = c(4.5, 4.5, 2.0, 2.0))
3621 plot(x, col = y + 9, pch = 19, cex = 1.5,
3622       xlim = c(0, 10), ylim = c(0, 10),
3623       xlab = bquote(x[1]), ylab = bquote(x[2]),
3624       cex.lab = 1.5, cex.axis = 1.5,
3625       main = 'SVM application to many points of two labels')
3626 axis(2, at = 0:10, tck = 1, lty = 3,
3627       col = "grey", labels = NA)
3628 axis(1, at = 0:10, tck = 1, lty = 3,
3629       col = "grey", labels = NA)
3630 lines(x1, x2p, lty = 2, col = 10)
3631 lines(x1, x2m, lty = 2, col = 11)
3632 lines(x1, x20, lwd = 1.5, col = 'purple')
3633 thetasvm = atan(-w[1]/w[2])*180/pi
3634 thetasvm
3635 #[1] 36.8699 #36.9 deg angle of the hyperplane
3636 #linear equations for the hyperplanes
3637 delx = 1.4
3638 dely = delx * (-w[1]/w[2])
3639 text(5 + 2*delx, 6.5 + 2*dely, bquote(bold(w%.%x) - b == 0),
3640       srt = thetasvm, cex = 1.5, col = 'purple')
3641 text(5 - delx, 7.6 - dely, bquote(bold(w%.%x) - b == 1),
3642       srt = thetasvm, cex = 1.5, col = 10)
3643 text(5, 4.8, bquote(bold(w%.%x) - b == -1),
3644       srt = thetasvm, cex = 1.5, col = 11)
3645 #normal direction of the hyperplanes
3646 arrows(2, 3.86, 2 + w[1], 4 + w[2], lwd = 2,
3647         angle = 15, length= 0.1, col = 'blue' )
3648 text(2 + w[1] + 0.4, 4 + w[2], bquote(bold(w)),
3649       srt = thetasvm, cex = 1.5, col = 'blue')
3650
3651 #new data points to be predicted
3652 xnew = matrix(c(0.5, 2.5, 7, 2, 6, 9),
3653                 ncol = 2, byrow = TRUE)
3654 points(xnew, pch = 17, cex = 2)
3655 predict(svmP, xnew) #Prediction
3656 #1 2 3
3657 #2 2 1
3658
3659 for(i in 1:3){
3660   text(xnew[i,1], xnew[i,2] - 0.4 ,
3661         paste('Q',i), cex = 1.5)
3662 }
3663 dev.off()

```

```

#Python for Fig. 9.7: SVM for a system of many points
# define x for plotting
x = np.array([1,6,2,8,3,7.5,1,8,
              4,9,5,9,3,7,5,9,1,5,
              5,3,6,4,7,4,8,6,9,5,
              10,6,5,0,6,5,8,2,2,2,1,1])
# define y
y = [1,1,1,1,1,1,1,1,1,2,2,2,2,2,2,2,2,2,2,2,2]
# redefine x for SVM
X = [[1,6],[2,8],[3,7.5],[1,8],
      [4,9],[5,9],[3,7],[5,9],[1,5],
      [5,3],[6,4],[7,4],[8,6],[9,5],
      [10,6],[5,0],[6,5],[8,2],[2,2],[1,1]]

svmP = svm.SVC(kernel='linear')
# train SVM
svmP.fit(X, y)

# predict new data using trained SVM
print('The support vectors are:', 
      svmP.support_vectors_, '\n')

# find hyperplane, normal vector, and SV (wx + b = 0)
w = -svmP.coef_[0]
print('w is:', w, '\n')

# rho is the negative intercept in R
# intercept is the positive intercept in Python
b = svmP.intercept_
print('b is:', b, '\n')

# find the maximum margin of separation
print(2/np.linalg.norm(w))
x1 = np.linspace(0,10,31)
x2 = (b - w[0]*x1)/w[1]
x2p = (1 + b - w[0]*x1)/w[1]
x2m = (-1 + b - w[0]*x1)/w[1]
x20 = (b - w[0]*x1)/w[1]

thetasvm = math.atan(-w[0]/w[1])*180/np.pi
# degree angle of the hyperplane
print(thetasvm)

delx = 1.4
dely = delx * (-w[0]/w[1])

# new predicted data points
svmP.predict([[0.5,2.5],[7,2],[6,9]])

```

```

#Python plot Fig. 9.7: SVM for a system of many points
#The plotting part of Fig. 9.7
plt.figure(figsize = (10,10))#set up plot
plt.xlim(-0.2,10.2)
plt.ylim(-0.2,10.2)
plt.title("SVM application to many points of two labels",
          pad = 10)
plt.xlabel("$x_1$", labelpad = 10)
plt.ylabel("$x_2$", labelpad = 10)

# plot points
color = ['r','r','r','r','r','r','r','r','forestgreen',
         'forestgreen','forestgreen','forestgreen',
         'forestgreen','forestgreen','forestgreen','forestgreen',
         'forestgreen','forestgreen','forestgreen']
plt.scatter(x[::2],x[1::2], color = color)
# plot newly predicted points
plt.scatter([0.5,7.6],[2.5,2.9], marker = "^",
            color = 'k', s = 90)
# plot lines
plt.plot(x1,x2p, color = 'red', linestyle = '--')
plt.plot(x1,x2m, color = 'forestgreen', linestyle = '--')
plt.plot(x1,x20,color = 'purple')

# add text
plt.text(5+2*delx,6.5+2*dely,'$w \cdot x - b = 0$',
         color = 'purple', rotation = thetasvm, size = 15)
plt.text(5-delx, 7.6-dely, '$w \cdot x - b = 1$',
         color = 'red', rotation = thetasvm, size = 15)
plt.text(5, 4.8,'$w \cdot x - b = -1$',
         color = 'forestgreen', rotation = thetasvm, size = 15)
plt.text(1.8,4.3,'$w$', color = 'blue', size = 15)
plt.text(0.3,2.1,'$Q_1$', color = 'k', size = 15)
plt.text(6.8,1.6,'$Q_2$', color = 'k', size = 15)
plt.text(5.8,8.6,'$Q_3$', color = 'k', size = 15)

# add normal direction of the hyperplanes
plt.arrow(2,3.86,w[0],w[1],color = 'blue', head_width = 0.1)
plt.show()

```

3665

3666 SVM maximizes $2/\|\mathbf{w}\|$, i.e., minimizes $\|\mathbf{w}\|$. The above computer code shows
 3667 that $\mathbf{w} = (-0.39996, 0.53328)$. This optimization is reached at the three SVs
 3668 $P_9(1,5), P_{17}(6,5)$ and $P_{19}(2,2)$. We can verify that the end points of these three
 3669 SVs are on the two hyperplanes, by plugging in these coordinates into the following
 3670 two linear equations

$$\mathbf{w} \cdot \mathbf{x} - b = \pm 1. \quad (3.45)$$

3671 Reversely, these equations have three parameters w_1, w_2, b , which can be determined
 3672 by three support vectors (SVs).

3673 The linear SVM separation can have nonlinear extensions, such as circular or
 3674 spherical hypersurfaces. Nonlinear SVM is apparently needed when the data cannot

3675 be easily separated by linear hyperplanes. Both R and Python codes have function
3676 parameters for the nonlinear SVM. Another SVM extension is from two classes to
3677 many classes. The multi-class SVM and nonlinear SVM are beyond the scope of
3678 this book. The interested readers are referred to the designated machine learning
3679 books, such as those listed at the section of References and Further Readings at
3680 the end of this chapter.

3681 **3.3 Random forest method for classification and**
3682 **regression**

3683 Random forest (RF) is a popular machine learning algorithm and belongs to the
3684 class of supervised learning. It uses many decision trees trained from the given data,
3685 called the training data. The training data usually include categorical data, such
3686 as weather types (e.g., rainy, cloudy, and sunny) as labels, and numeric data, such
3687 as historical instrument observations (e.g., atmospheric pressure, humidity, wind
3688 speed, wind direction, and air temperature). These data form many logical decision
3689 trees that decide under what sets of conditions the weather is rainy, cloudy, or
3690 sunny. Many decision trees form a forest. Multiple decision nodes and branches for
3691 each tree are determined and trained by using a set of random sampling procedures.
3692 The random sampling and the decision trees lead to the name of random forest.
3693 The training step results in an RF model which are a set of decision trees. In the
3694 prediction step, you submit your new data to the the trained RF model whose each
3695 decision tree makes a vote for a category based on your new data. If the category
3696 “tomorrow is rainy” receives the most votes, then you have predicted a rainy day
3697 tomorrow. In this sense, RF is an ensemble learning method. Because of this nature
3698 of ensemble predictions, an RF algorithm should ensure the votes among the trained
3699 trees should have as little correlation as possible.
3700

3701 The above is a layman’s description of RF for classification, when the RF objective
3702 is for categorical data, or called factor in data types. RF can also be used to fill
3703 in the missing numeric real values. The missing values, or missing data in climate
3704 science, are often denoted by NA, NaN, or -99999. RF makes a prediction for the
3705 missing data. The prediction is an ensemble mean from the trained RF trees. This
3706 is an RF regression, which produces numerical results.

3707 **3.3.1 RF flower classification for a benchmark iris dataset**
3708

3709 We use the popular iris flower dataset (Fisher 1936) frequently used for the RF
3710 teaching as an example to explain the RF computing procedures and to interpret
3711 the RF results. Sir Ronald A. Fisher (1890 – 1962) was a British statistician and
3712 biologist. The Fisher dataset contains three iris species: *setosa*, *versicolor*, *virginica*.
3713 The data are the length and width of sepal and petal of a flower. The first two rows
3714 of the total 150 rows of the dataset are below.

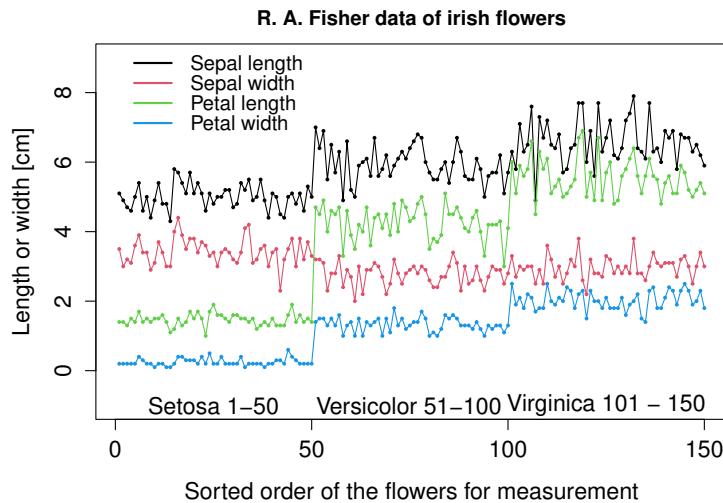


Fig. 3.8 The R. A. Fisher iris dataset of sepal length, sepal width, petal length, and petal width for flower species (Fisher 1936).

```

3715 #   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
3716 #1      5.1        3.5        1.4        0.2  setosa
3717 #2      4.9        3.0        1.4        0.2  setosa

```

3718 Figure 3.8 shows the entire dataset.

3719 Figure 3.8 may be generated by the following computer code.

```

3720 #R plot Fig. 9.8: R.A. Fisher data of three iris species
3721 setwd('~/climstats')
3722 data(iris) #read the data already embedded in R
3723 dim(iris)
3724 #[1] 150 5
3725 iris[1:2,] # Check the first two rows of the data
3726 #   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
3727 #1      5.1        3.5        1.4        0.2  setosa
3728 #2      4.9        3.0        1.4        0.2  setosa
3729
3730 str(iris) # Check the structure of the data
3731 #'data.frame': 150 obs. of 5 variables:
3732 #$ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
3733 #$ Species     : Factor w/ 3 levels "setosa","versicolor",...
3734
3735 setEPS() #Plot the data of 150 observations
3736 postscript("fig0908.eps", width=7, height=5)
3737 par(mar= c(4.5, 4.5, 2.5, 0.2))
3738 plot(iris[,1], type = 'o', pch = 16, cex = 0.5, ylim = c(-1, 9),
3739       xlab = 'Sorted_order_of_the_flowers_for_measurement',
3740       ylab = 'Length_or_width_[cm]',
3741       main = 'R.A.Fisher_data_of_irish_flowers', col = 1,
3742       cex.lab = 1.3, cex.axis = 1.3)

```

```

3743 lines(iris[,2], type = 'o', pch = 16, cex = 0.5, col=2)
3744 lines(iris[,3], type = 'o', pch = 16, cex = 0.5, col=3)
3745 lines(iris[,4], type = 'o', pch = 16, cex = 0.5, col=4)
3746 legend(0, 9.5, legend = c('Sepal.length', 'Sepal.width',
3747           'Petal.length', 'Petal.width'),
3748           col = 1:4, lty = 1, lwd = 2, bty = 'n',
3749           y.intersp = 0.8, cex = 1.2)
3750 text(25, -1, 'Setosa_1-50', cex = 1.3)
3751 text(75, -1, 'Versicolor_51-100', cex = 1.3)
3752 text(125, -1, 'Virginica_101-150', cex = 1.3)
3753 dev.off()

```

```

# Python plot Fig. 9.8: R.A. Fisher data of iris species
iris = datasets.load_iris()
# notice "target" is equivalent to "species"
iris = pd.DataFrame(data=np.c_[iris['data'],iris['target']],
                      columns=iris['feature_names'] + ['target'])

# check the structure of the data
print(iris.info(), '\n')

# check the first two rows of the data
iris.iloc[0:2,:]

# set up plot
plt.title("R.A. Fisher data of iris flowers", pad = 10)
plt.xlabel("Sorted order of the flowers for measurement",
           labelpad = 10)
plt.ylabel("Length or width [cm]", labelpad = 10)
plt.xticks([0,50,100,150])
plt.ylim(-2,9)
plt.yticks([0,2,4,6,8])

# plot the data of 150 observations
x = np.linspace(0,149,150)
plt.plot(x,iris.iloc[:,0],'ko-', label = 'Sepal.length')
plt.plot(x,iris.iloc[:,1],'ro-', label = 'Sepal.width')
plt.plot(x,iris.iloc[:,2],'go-', label = 'Petal.length')
plt.plot(x,iris.iloc[:,3],'bo-', label = 'Petal.width')

# add legend
plt.legend()

# add text
plt.text(13,-1, "Setosa_1-50", size = 15)
plt.text(57,-1, "Versicolor_51-100", size = 15)
plt.text(107,-1, "Virginica_101-150", size = 15)

plt.show()

```

3754

The first 50 are setosa, which are the smallest flower, and the last 50 for virginica that have the largest petals. Given the data of sepal and petal sizes, one may correctly detect their corresponding flower species. However, not all the flowers have

3758 the same size. This makes the species detection more difficult. The random forest
 3759 algorithm can make the detection with a small error. We make an RF experiment
 3760 with this dataset. We randomly select 120 rows of data as the training data to
 3761 obtained a trained RF model with 500 decision trees, use the remained 30 rows
 3762 as the new data for our detection. These decision trees to vote on each of the 30
 3763 rows of the new data. The most species votes a row receives from the trained RF
 3764 trees is the RF detected result. By default, the decision trees in an RF algorithm
 3765 are built randomly, the RF result from each run based on the same data may be
 3766 slightly different. The following computer code shows the RF run and its result of
 3767 this experiment for Fisher's iris data. The code also generates Fig. 3.9.

```

3768 #R code: RF prediction using the Fisher iris data
3769 #install.packages("randomForest")
3770 library(randomForest)
3771 #set.seed(8) # run this line to get the same result
3772 #randomly select 120 observations as training data
3773 train_id = sort(sample(1:150, 120, replace = FALSE))
3774 train_data = iris[train_id, ]
3775 dim(train_data)
3776 #[1] 120 5
3777 #use the remaining 30 as the new data for prediction
3778 new_data = iris[-train_id, ]
3779 dim(new_data)
3780 #[1] 30 5
3781
3782 #train the RF model
3783 classifyRF = randomForest(x = train_data[, 1:4],
3784                           y = train_data[, 5], ntree = 800)
3785 classifyRF #output RF training result
3786 #Type of random forest: classification
3787 #Number of trees: 800
3788 #No. of variables tried at each split: 2
3789
3790 #OOB estimate of error rate: 4.17%
3791 #Confusion matrix:
3792 #
3793 #          setosa versicolor virginica class.error
3794 #setosa      41         0         0  0.00000000
3795 #versicolor   0        34         2  0.05555556
3796 #virginica    0         3        40  0.06976744
3797
3798 plot(classifyRF) #plot the errors vs RF trees Fig. 9.9a
3799
3800 classifyRF$importance #classifyRF$ has many outputs
3801 #          MeanDecreaseGini
3802 #Sepal.Length     8.313131
3803 #Sepal.Width      1.507188
3804 #Petal.Length     31.075960
3805 #Petal.Width      38.169763
3806 varImpPlot(classifyRF) #plot the importance result Fig. 9.9b
3807
3808 #RF prediction for the new data based on the trained trees
3809 predict(classifyRF, newdata = new_data[,1:4])
3810 # 2           4           10          11          12          14
3811 #setosa       setosa       setosa       setosa       setosa       setosa

```

```

3811 #It got two wrong: 120 versicolor and 135 versicolor
3812
3813
3814 #Another version of the randomForest() command
3815 anotherRF = randomForest(Species ~ .,
3816                               data = train_data, ntree = 500)

```

```

#Python code: RF prediction with the Fisher iris data
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn import metrics
# randomly select 120 observations as training data
train_id = random.sample(range(0, 150), 120)
train_data = iris.iloc[train_id,:]
print(train_data.shape)# dimensions of training data
# use the remaining 30 as the new data for prediction
new_data = iris.drop(train_id)
print(new_data.shape)# dimensions of new data
# train the RF model
classifyRF = RandomForestClassifier(n_estimators=800,
                                      oob_score=True)
training = classifyRF.fit(train_data.iloc[:,4],
                           train_data.iloc[:,4])
predictions = classifyRF.predict(new_data.iloc[:,4])
print(predictions)
#[0. 0. 0. 0. 0. 0. 0. 0. 1. 1. 1. 2.
#1. 2. 1. 2. 1. 1. 2. 2. 2. 2.
# 2. 2. 2. 2. 2.]
#0 = setosa, 1 = versicolor, 2 = virginica
confusion_M = confusion_matrix(new_data.iloc[:,4],
                                 predictions)
print(confusion_M)
#[[ 9  0  0] #all correct predictions for setosa
#[ 0  8  3] #3 wrong predictions for versicolor
#[ 0  0 10]] #all correct predictions for virginica
accuracy_score(new_data.iloc[:,4],
                predictions)
#0.9

```

```
#Python lot Fig. 9.9(a): RF mean accuracy
n_estimators = 100
classifyRF = RandomForestClassifier(n_estimators,
                                     oob_score=True)
RFscore = []
for i in range(1, n_estimators + 1):
    classifyRF.set_params(n_estimators=i)
    classifyRF.fit(train_data.iloc[:, :4],
                   train_data.iloc[:, 4])
    RFscore.append(classifyRF.score(train_data.iloc[:, :4],
                                    train_data.iloc[:, 4]))
plt.plot(RFscore)
plt.title("Random forest score of the mean accuracy")
plt.xlabel("Number of estimators")
plt.ylabel("OOB score")
plt.show()
```

3818

```
#Python plot Fig. 9.9(b): RF importance plot
# Python shows the importance results differently from R
plt.plot(np.linspace(1, 4, 4),
         classifyRF.feature_importances_,
         'ko', markersize = 15)
plt.title("Importance plot of RF model", pad = 10)
plt.ylabel("Mean decrease in impurity", labelpad = 12)
plt.xticks([1, 2, 3, 4], ['Sepal.Length', 'Sepal.Width',
                        'Petal.Length', 'Petal.Width'], rotation = 90)
plt.yticks(classifyRF.feature_importances_,
           [0.01, 0.11, 0.39, 0.48])
plt.grid()
plt.show()
```

3819

3820 The following explains the RF output data and figures.

- 3821 • Confusion matrix: This matrix displays correct and wrong classifications based
3822 on the training data. The columns are truth and rows are RF classifications.
3823 The diagonals are corrected classifications, and off diagonals are the wrong
3824 ones. The 41 setosas are all correctly classified. Among the 37 versicolors, 34
3825 are correctly classified, but 3 are wrongly classified as virginica. Among the 37
3826 versicolors, 34 are correctly classified, but 3 are wrongly classified as virginica.
- 3827 • Classification error: The last column in the confusion matrix is the classification
3828 error which is equal to the sum of the off-diagonal elements in that row divided
3829 by the sum of the entire row, i.e., the ratio of the number of wrong classifications
3830 to the total number of classifications. The errors shows how good is the trained
3831 RF model.
- 3832 • OOB estimate of error rate: OOB stands for out-of-bag, a term for a statistical
3833 sampling process that puts a part of the data for training the RF model and
3834 the remainder for validation. This remainder is referred to as the OOB set

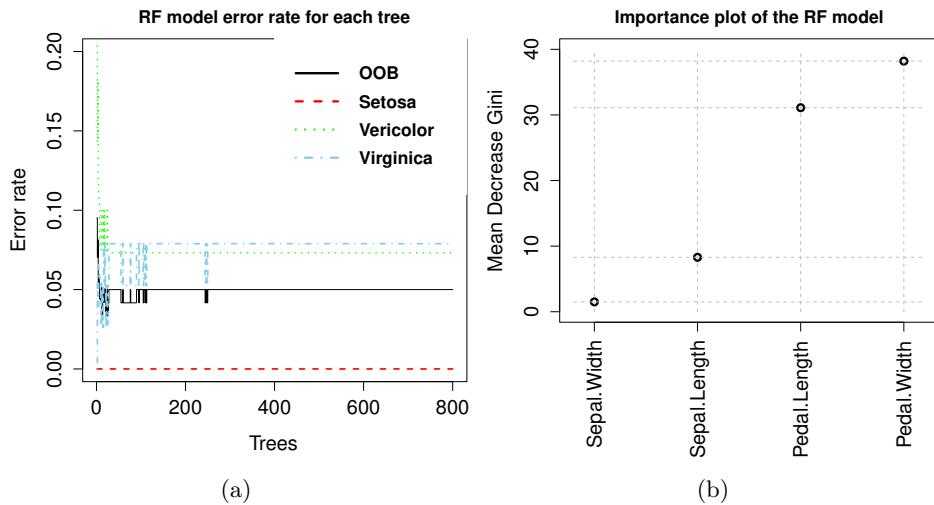


Fig. 3.9 (a) RF errors vs. trees. (b) The RF's importance plot.

of data. The OOB error rate is defined as the ratio of the number of wrong decisions based on the majority vote to the size of OOB set.

- No. of variables tried at each split: This is the number of variables randomly sampled for growing trees, and is denoted by `mtry`. The Fisher iris data have four variables ($p = 4$). It is recommended that $mtry = \sqrt{p}$ for RF classification, and $= p/3$ for RF regression.
- Number of trees: We want to build enough trees for RF so that the RF errors shown in Fig. 3.9(a) is stabilized. Too few trees may yield results of large differences at different runs.
- RF errors vs. trees: The colored lines show errors of different types of validation procedures. As RF users, we pay attention to whether these errors are stabilized. The detailed error definitions are beyond the scope of this chapter.
- Mean decrease Gini scores and the importance plot (Fig. 3.9(b)): A higher mean decrease Gini (MDG) score corresponds to more importance of the variable in the model. In our example, the petal width is the most important variable. This agrees with our intuition from Fig. 3.8. The figure shows that the petal width has clear distinctions among the three species and have little variance. The petal length also has clear distinctions, but has larger variances, and is the second most important variable. The sepal width has almost no distinctions among the species and has the smallest MDG score 1.5. It is the least important variable.
- RF prediction: Finally, the RF predictions were made for the 30 rows of new data. The prediction results for the first six rows of the new data are shown here. These are the final results. Among the 30 rows, RF correctly predicted 28 and got only two wrong: 120 and 135 should be virginica, but RF predicted versicolor for both. In weather forecast, these can be the public weather outlook

3861 of the 9th day from today (e.g., sunny, rainy, or cloudy) for 30 different locations
 3862 in a country.

3.3.2 RF regression for the daily ozone data of New York City

3863 Similar to Fisher's iris data, the daily air quality data, the file name `airquality`,
 3864 of New York City (NYC) from 1 May 1973 - 30 September 1973 (153 days) is an
 3865 other benchmark data for the random forest algorithms. The dataset contains the
 3866 following weather parameters: Ozone concentration in the ground-level air, cumu-
 3867 lative solar radiation, average wind speed, and daily maximum air temperature.
 3868 When the ozone concentration is higher than 86 parts per billion (ppb), the air is
 3869 considered unhealthy. The ozone concentration is related to solar radiation, wind
 3870 and temperature. The following list provides more information on this `airquality`
 3871 dataset:

- 3872 • The ozone data in ppb observed between 1300 and 1500 hours at Roosevelt Island
 of the New York City. Among the 153 days, 37 had no data and are denoted by
 NA. The data range is 1 - 168 ppb. The maximum ozone level of this dataset
 is 168 ppb, occurred on 25 August 1973, when the cumulative solar radiation
 was high at 238 Lang, the average wind speed low at 3.4 mph, and Tmax is
 moderate 81°F.
- 3873 • Cumulative solar radiation from 0800 to 1200 hours with units in Langley (Lang
 or Ly) (1 Langley = 41, 868 Watt·sec/m², or 1 Watt/m² = 0.085985 Lang-
 ley/hour) in the lightwave length range 4,000 - 7,700 Angstroms at Central
 Park of the New York City. The data range is 7 - 334 Lang.
- 3874 • Average wind speed in miles per hour (mph) at 0700 and 1000 hours at LaGuardia
 Airport, less than 10 km from Central Park. The data range is 1.7 - 20.7 mph.
- 3875 • Maximum daily temperature Tmax in °F at La Guardia Airport. The data
 range is 56 - 97 °F. The daily Tmax data can also be downloaded from the
 NOAA NCEI website by an online search using the words: NOAA Climate Data
 Online LaGuardia Airport.

3876 While RF analysis for this `airquality` dataset can be for many purposes, our
 3877 example is to fill the 37 missing ozone data with RF regression results. RF use the
 3878 above four parameters and their data to build decision trees. Figure 3.10(a) shows
 3879 the original 116 observed data points in black dots and lines, and the RF regression
 3880 estimates for the 37 missing data (the blue dots and lines). A single blue dot means
 3881 only an isolated day of missing data. A blue line indicates missing data in successive
 3882 days. For better visualization, Fig. 3.10(b) shows the complete NYC daily ozone
 3883 time series that join the observed data with the result data from an RF regression.

3884 Figure 3.10 may be generated by the following computer code.

```
3885 #R plot Fig. 9.10: RF regression for ozone data
3886 library(randomForest)
3887 airquality[1:2,] #use R's RF benchmark data "airquality"
3888 # Ozone Solar.R Wind Temp Month Day
```

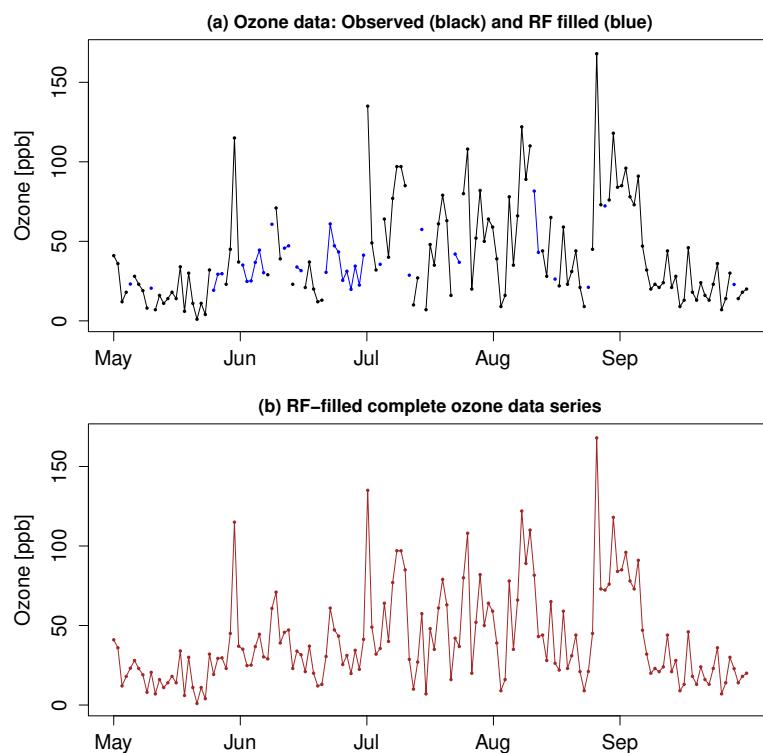


Fig. 3.10 (a) New York City ozone data from 1 May to 30 September 1973. Among the 153 days, 116 had observed data (black dots and lines) and 37 had missing data. The missing data are filled by RF regression results (blue dots and lines). (b) The complete ozone data time series as a continuous curve with dots when the missing data are replaced by the RF regression.

```

3903 #1      41      190   7.4    67      5    1
3904 #2      36      118   8.0    72      5    2
3905 dim(airquality)
3906 #[1] 153   6
3907 ozoneRFreg = randomForest(Ozone ~ ., data = airquality,
3908           mtry = 2, ntree = 500, importance = TRUE,
3909           na.action = na.roughfix)
3910 #na.roughfix allows NA to be replaced by medians
3911 #to begin with when training the RF trees
3912 ozonePred = ozoneRFreg$predicted #RF regression result
3913 t0 = 1:153
3914 n1 = which(airquality$Ozone > 0) #positions of data
3915 n0 = t0[-n1] #positions of missing data
3916 ozone_complete = ozone_filled= airquality$Ozone
3917 ozone_complete[n0] = ozonePred[n0] #filled by RF
3918 ozone_filled = ozonePred #contains the RF reg result
3919 ozone_filled[n1] <- NA #replace the n1 positions by NA

```

```
3920 t1 = seq(5, 10, len = 153) #determine the time May - Sept
3921
3922 par(mfrow = c(2, 1))
3923 par(mar = c(3, 4.5, 2, 0.1))
3924 plot(t1, airquality$Ozone,
3925     type = 'o', pch = 16, cex = 0.5, ylim = c(0, 170),
3926     xlab = '', ylab = 'Ozone[ppb]', xaxt="n",
3927     main = '(a) Ozone data: Observed(black) and RF filled(blue)',
3928     col = 1, cex.lab = 1.3, cex.axis = 1.3)
3929 MaySept = c("May", "Jun", "Jul", "Aug", "Sep")
3930 axis(side=1, at=5:9, labels = MaySept, cex.axis = 1.3)
3931 points(t1, ozone_filled, col = 'blue',
3932         type = 'o', pch = 16, cex = 0.5)#RF filled data
3933
3934 #Plot the complete data
3935 par(mar = c(3, 4.5, 2, 0.1))
3936 plot(t1, ozone_complete,
3937     type = 'o', pch = 16, cex = 0.5, ylim = c(0, 170),
3938     xlab = '', ylab = 'Ozone[ppb]',
3939     xaxt="n", col = 'brown',
3940     main = '(b) RF-filled complete ozone data series',
3941     cex.lab = 1.3, cex.axis = 1.3)
3942 MaySept = c("May", "Jun", "Jul", "Aug", "Sep")
3943 axis(side=1, at=5:9, labels = MaySept, cex.axis = 1.3)
```

```

# Python plot Fig. 9.10: RF regression for ozone data
from sklearn.ensemble import RandomForestRegressor
# Read in the airquality data and create a copy
airquality = pd.read_csv("data/airquality.csv",
                         index_col = 'Unnamed: 0')
airquality_copy = airquality.copy()
print(airquality.head())#print the first 5 rows of data
print(np.shape(airquality))#data matrix dim (153, 6)
# Create list of integer from 1 to 153
t0 = list(np.linspace(1,153,153, dtype=int))
# positions of recorded data
n1 = np.where(airquality["Ozone"] > 0)
n1 = n1[0].tolist()
n1 = [x+1 for x in n1]
# positions of unknown 'NaN' data
n0=[]
for x in t0:
    if x not in n1:
        n0.append(x)
# Replace NA with medians in order to train the RF trees
airquality['Solar.R'].fillna(
    value=airquality['Solar.R'].median(), inplace=True)
airquality['Wind'].fillna(
    value=airquality['Wind'].median(), inplace=True)
airquality['Temp'].fillna(
    value=airquality['Temp'].median(), inplace=True)
airquality['Day'].fillna(
    value=airquality['Day'].median(), inplace=True)
airquality['Month'].fillna(
    value=airquality['Month'].median(), inplace=True)
airquality['Ozone'].fillna(
    value=airquality['Ozone'].median(), inplace=True)
# Create our features X and our target y
X = airquality[['Solar.R', 'Wind', 'Temp', 'Month', 'Day']]
y = airquality[['Ozone']]
# split our data into training and test sets
X_train = X.loc[n1,:]
X_test = X.loc[n0,:]
y_train = y.loc[n1,:]
y_test = y.loc[n0,:]
#create the RF Regressor model with 500 estimators
ozoneRFreg = RandomForestRegressor(n_estimators=500,
                                    oob_score=True, max_features = 2)
# fit model to our training set
ozoneRFreg.fit(X_train, y_train.values.ravel())
#create our prediction
ozonePrediction = ozoneRFreg.predict(X_test)
#create data frame with our results
result = X_test
result['Ozone'] = y_test
result['Prediction'] = ozonePrediction.tolist()
#create data frame with our predictions
ozone_filled = pd.DataFrame(None, index = np.arange(153),
                            columns = ["Prediction"])
ozone_filled.loc[n0, 'Prediction'] = result["Prediction"]
ozone_filled.index += 1

```

```
#Python plot Fig. 9.10(a)
t1 = np.linspace(5,10,153)
mfrow = [2,1]
mar = [3, 4.5, 2, 0.1]
# original known data
ar = pd.read_csv("data/airquality.csv",
                  index_col = 'Unnamed:0')
plt.plot(t1, ar["Ozone"],
          'ko-', markersize = 3)
plt.ylim(0, 170)
#unknown Predicted data
plt.plot(t1, ozone_filled['Prediction'],
          'bo-', markersize = 3)
plt.ylabel("Ozone[ppb]", labelpad = 12)
MaySept = ["May", "Jun", "Jul", "Aug", "Sep"]
plt.xticks([5,6,7,8,9], MaySept)
plt.title("(a) Observed(black) and RF filled(blue)")
plt.show()
```

3945

```
#Python plot Fig. 9.10(b)
# combine unknown and known data into one dataframe
result_copy = result.copy()
result_copy = result_copy.rename(
    columns = {'Ozone':'y_tested', 'Prediction':'Ozone'})
ozone_complete = airquality_copy[
    "Ozone"].fillna(result_copy["Ozone"])
#plot the combined data
t1 = np.linspace(5,10,153)
mfrow = [2,1]
mar = [3, 4.5, 2, 0.1]
plt.plot(t1, ozone_complete, 'ro-', markersize = 3)
plt.ylim(0, 170)
plt.ylabel("Ozone[ppb]", labelpad = 12)
MaySept = ["May", "Jun", "Jul", "Aug", "Sep"]
plt.xticks([5,6,7,8,9], MaySept)
plt.title("(b) RF-filled complete ozone data series")
plt.show()
```

3946

3.3.3 What does a decision tree look like?

We have learned that the RF prediction is the result of majority votes from the trained decision trees in a random forest. Then, what does a decision tree look like? Figure 3.11 shows a decision tree in the RF computing process for the training data of iris flowers. The gray box on top shows the percentage of each species in the 120 rows of training data. This particular set of training data was randomly selected from the entire R.A. Fisher dataset which has 150 rows. Among the 120 rows of training data, 41 rows are setosa (34%), 42 rows versicolor (35%), and 37

3956 rows (31%) virginica. The first branch of the decision tree grows from a condition

$$\text{Petal Length} < 2.5 \text{ [cm]}. \quad (3.46)$$

3957 If “yes”, this iris is setosa. All the 41 setosa flowers (34% of the entire training
 3958 data) have been detected. This decision is clearly supported by the real petal length
 3959 data (the green line) shown in Fig. 3.8, which shows that only setosa has petal length
 3960 less than 2.5.

3961 If “no”, then two possibilities exist. This allows us to grow a new branch of the
 3962 tree. The condition for this branch is

$$\text{Petal Width} < 1.8 \text{ [cm]}. \quad (3.47)$$

3963 This condition determines whether a flower is versicolor or virginica. After the
 3964 isolation of 41 setosa flowers, the remaining 79 flowers are 42 versicolor (53%)
 3965 and 37 virginica (47%). The “yes” result of condition (3.47) leads to versicolor.
 3966 This result is 93% correct and 7% incorrect. The “no” result implies virginica. This
 3967 conclusion is 97% correct, and 3% incorrect. These conclusions are supported by the
 3968 real petal width data (the blue line) shown in Fig. 3.8. The petal width of versicolor
 3969 iris flowers is in general less than 1.8 [cm]. However, the data have some fluctuations
 3970 in both the versicolor and virginica sections in Fig. 3.8. These fluctuations lead to
 3971 the errors of decision.

3972 The entire RF process for our iris species example had grown 800 such decision
 3973 trees using the 120 rows of training data. Different trees use different branch
 3974 conditions. RF algorithms grow these trees following various kinds of optimization
 3975 principles, which involve some tedious mathematics not covered here.

3976 Figure 3.11 may be generated by the following computer code.

```
3977 #R plot of Fig. 9.11: A tree in a random forest
3978 #install.packages('rpart')
3979 library(rpart)
3980 #install.packages('rpart.plot')
3981 library(rpart.plot)
3982
3983 setwd('/~/climstats')
3984 setEPS() #Plot the data of 150 observations
3985 postscript("fig0911.eps", width=6, height=6)
3986 par(mar = c(0, 2, 1, 1))
3987 iris_tree = rpart( Species ~ . , data = train_data)
3988 rpart.plot(iris_tree,
3989   main = 'An decision tree for the RF training data')
3990 dev.off()
```

An decision tree for the RF training data

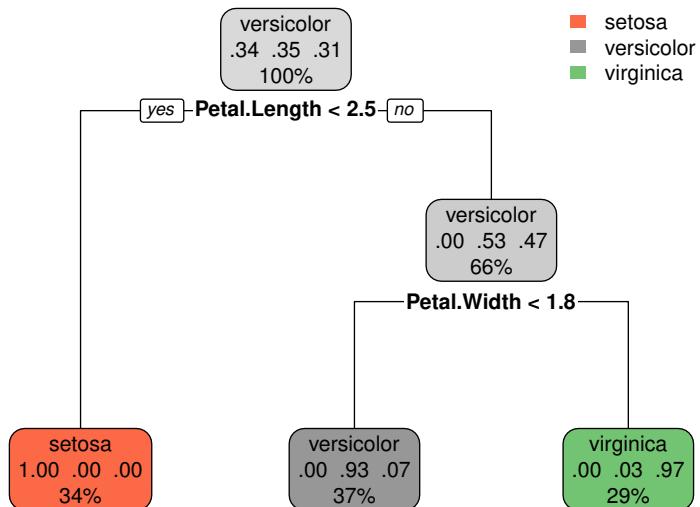


Fig. 3.11 A trained tree in the random forest for the iris data of R.A. Fisher (1936).

```
# Python plot Fig. 9.11: A tree in a random forest
classifyRF = RandomForestClassifier(n_estimators=1,
                                     oob_score=True)
training = classifyRF.fit(train_data.iloc[:, :4],
                           train_data.iloc[:, 4])
features = ['sepal_length(cm)', 'sepal_width(cm)',
            'petal_length(cm)', 'petal_width(cm)']
class_name = ['setosa', 'versicolor', 'virginica']
# set up figure
plt.figure(figsize = (15,15))
plot_tree(classifyRF.estimators_[0],
          feature_names = features,
          class_names = class_name, filled = True,
          fontsize = 12)
plt.show()
```

3992

3.4 Neural network and deep learning

3993

3994 A neural network (NN) consists of a series of data fitting based on the training data
3995 and an application of the fitted NN model for test data. The NN outputs categorical
3996 predictions, such as a sunny or rainy day, or numerical predictions as a regression.
3997 NN is also known as an artificial neural network (ANN) or simulated neural network
3998 (SNN). It is a popular machine learning method, and is a fundamental building
3999 block the deep learning algorithms that often involve the data fitting of multiple
4000 layers, referred to as hidden layers.

4001 Why is it called neural network? How does a data fitting process have anything
4002 to do with “neural” and/or “network”? Artificial neurons were first proposed by
4003 Warren Sturgis McCulloch (1898 - 1969), an American neurophysiologist, and Wal-
4004 ter Harry Pitts (1923 - 1969), an American logician, in their 1943 paper entitled “
4005 A logical calculus of ideas immanent in nervous activity.” This mathematical paper
4006 used terms “neuron”, “action”, “logic expression”, and “net”, which are among
4007 the keywords in the modern NN writings. The ten theorems of the paper formu-
4008 late a suite of logic expressions based on data. The word “neuron” was more a
4009 graphic indication for actions and logic expressions than a biological reference. A
4010 layman’s understanding of an NN machine learning is often put incorrectly toward
4011 biological neurons and a biological neural network. Therefore, ANN may be a more
4012 appropriate term for NN to avoid confusion.

4013 This book attempts to make such a brief NN introduction that you can use
4014 our R or Python code for your data and objectives, interpret your NN computing
4015 results, and understand the basic principles of mathematical formulations of an NN
4016 algorithm. However, we do not attempt to derive mathematical details of the NN
4017 theory. This section has two sub-sections: (i) A simple NN example of a decision
4018 system, and (ii) An example of NN prediction using the benchmark data of Fisher’s
4019 iris flowers.

4020

4021

3.4.1 An NN model for an automated decision system

4022

4023

3.4.1.1 An overall idea of a neural network decision system for recruitment

4024 A senior human resource manager of an IT company recruited three new employ-
4025 ees among six candidates. Her hiring decisions were made based on the technical
4026 knowledge scores (TCS) and communication skills scores (CSS) given in Table 9.1.
4027 TCS and CSS were the interview results from company’s working groups. A ju-
4028 nior human resource manager wishes to use an NN model to help him make his
4029 recruitment decisions consistent with those made by the senior manager.

4030 The junior manager has received the following scores for the three new job appli-
4031 cants A, B, and C, whose TKS and CSS scores are as follows: TKS (30, 51, 72), and
4032 CSS (85, 51, 30). The senior manager’s ruling seems to suggest that a candidate is

Table 9.1 TCS and CSS data and recruitment decisions

| | | | | | | |
|----------|------|--------|--------|--------|------|------|
| TKS | 20 | 10 | 30 | 20 | 80 | 30 |
| CSS | 90 | 20 | 40 | 50 | 50 | 80 |
| Decision | Hire | Reject | Reject | Reject | Hire | Hire |

4033 recruited when either TKS score or CSS score is high. Thus, this junior manager
 4034 has an easy decision for Candidate A whose CSS 85 is high, so to be recruited.
 4035 However, the decision for Candidate B is difficult. Neither TKS 51 nor CSS 51 of
 4036 Candidate B is high, but both are higher than those of the rejected, and further the
 4037 sum of TKS and CSS is 102, lower than the minimum of the totals of the recruited,
 4038 110. Should the junior manager hire Candidate B? The decision for Candidate C is
 4039 also difficult. Candidate C seems weaker than the recruited case TKS 30 and CSS
 4040 80, but not too much weaker. Should the junior manager hire Candidate C? The
 4041 following NN computer code may learn from the data of the senior manager and
 4042 suggest an NN decision that may serve as a useful reference for the junior manager.

4043 **3.4.1.2 An NN code, results, and their interpretation**

```

4044 #R code for NN recruitment decision and Fig. 9.12
4045 #Ref: https://www.datacamp.com/tutorial/neural-network-models-r
4046 TKS = c(20,10,30,20,80,30)
4047 CSS = c(90,20,40,50,50,80)
4048 Recruited = c(1,0,0,0,1,1)
4049 #make a data frame for the NN function neuralnet
4050 df = data.frame(TKS, CSS, Recruited)
4051
4052 require(neuralnet) # load 'neuralnet' library
4053 # fit neural network
4054 set.seed(123)
4055 nn = neuralnet(Recruited ~ TKS + CSS, data = df,
4056                 hidden = 5, act.fct = "logistic",
4057                 linear.output = FALSE)
4058 plot(nn) #Plot Fig 9.12: A neural network
4059
4060 TKS=c(30,51,72) #new data for decision
4061 CSS=c(85,51,30) #new data for decision
4062 test=data.frame(TKS,CSS)
4063 Predict=neuralnet::compute(nn,test)
4064 Predict$net.result #the result is probability
4065 #[1,] 0.99014936
4066 #[2,] 0.58160633
4067 #[3,] 0.01309036
4068
4069 # Converting probabilities into decisions
4070 ifelse(Predict$net.result > 0.5, 1, 0) #threshold = 0.5
4071 #[1,] 1
4072 #[2,] 1
4073 #[3,] 0
4074
4075 #print bias and weights
4076 nn$weights[[1]][[1]]

```

```
4077 #print the last bias and weights before decision
4078 nn$weights[[1]][[2]]
4079 #print the random start bias and weights
4080 nn$startweights
4081 #print error and other technical indices of the nn run
4082 nn$result.matrix #results data
```

```

#Python code for the NN recruitment decision and Fig. 9.12
#imports to build neural network model and visualize
#You need to install the following Python environments:
#tensorflow and keras-models
#!pip3 install ann_visualizer
#!pip3 install keras-models
#from a terminal: pip install tensorflow
from ann_visualizer.visualize import ann_viz
#from graphviz import Source
#from tensorflow.keras import layers
from keras.models import Sequential
from keras.layers import Dense

TKS = [20,10,30,20,80,30]
CSS = [90,20,40,50,50,80]
Recruited = [1,0,0,0,1,1]
#combine multiple columns into a single set
df = pd.DataFrame({'TKS': TKS,'CSS': CSS,
                    'Recruited': Recruited})
df.head()
X = df[['TKS', "CSS"]]
y = df[['Recruited']]
#random.seed(123)
#The model is random due to too little training data

nn = Sequential()
nn.add(Dense(5, input_dim = 2, activation = "sigmoid"))
nn.add(Dense(1,activation = "sigmoid"))
#nn.add(Dense(2,activation = "sigmoid"))
nn.compile(optimizer = "adam", loss = "BinaryCrossentropy",
            metrics = "BinaryAccuracy")
#fit neural network to data
nn.fit(X,y)
#visualize the neural network
ann_viz(nn, title = "Recruitment_Decision")
#output the nn model weights and biases
print(nn.get_weights())

TKS = [30,51,72] #new data for decision
CSS = [85,51,30] #new data for decision
test = pd.DataFrame({'TKS': TKS,'CSS': CSS})
prediction = nn.predict(test)
print(prediction)
#[[0.70047873]
#[0.66671777]
#[0.38699177]]
#Converting probabilities into decisions
for i in range(3):
    if prediction[i] >= 0.5:
        print("Hire")
    else:
        print("Reject")
#Hire
#Hire
#Reject

```

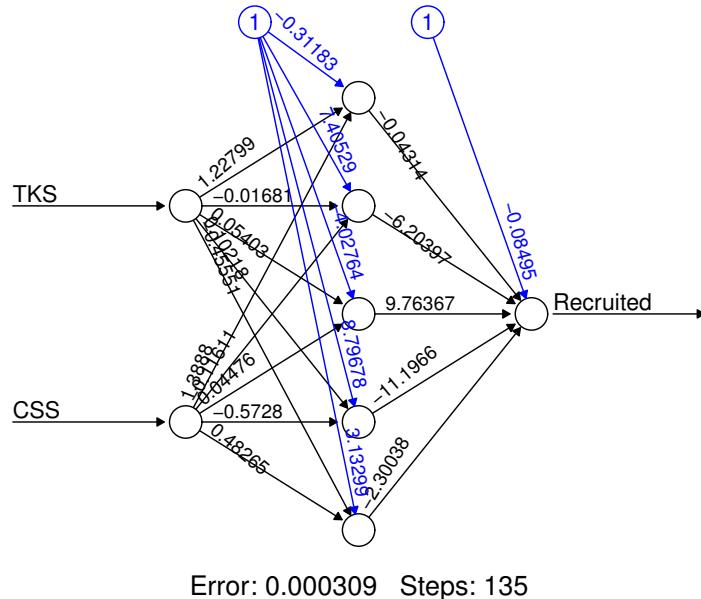


Fig. 3.12 A simple neural network of five neurons in a hidden layer for an NN hiring decision system.

Figure 3.12 shows a simple neural network plotted by the above computer code. The black numbers are called weights w_{ij} that are multiplied by the input data x_i for neuron j . The blue numbers are called biases b_j , associated with neurons j , indicated by the circles pointed by a blue arrow. The last circle on the right is the result, or called output layer. The first two circles on the left indicate input data, and form the input layer. The weights and biases are mathematically aggregated in the following way:

$$z_j = \sum_{i=1}^n w_{ij}x_i + b_j, \quad (3.48)$$

where z_j are for fitting an activation function at neuron j when all the training data are used. In the above code, the activation function is logistic. A logistic function is defined as

$$g(z) = \frac{1}{1 + \exp(-k(z - z_0))}, \quad (3.49)$$

where k is called the logistic growth rate, and z_0 is the midpoint. This function can also be written as

$$g(z) = \frac{1}{2} + \frac{1}{2} \tanh\left(-\frac{k(z - z_0)}{2}\right), \quad (3.50)$$

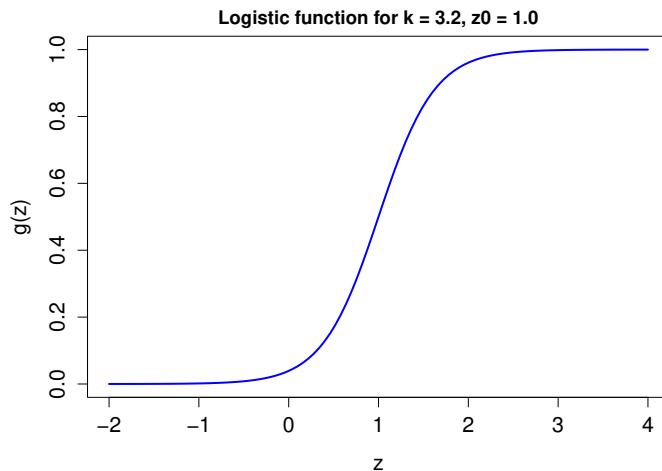


Fig. 3.13 Logistic function with growth rate $k = 3.2$ and midpoint $y_0 = 1.0$.

4096 where the tanh function is defined as

$$\tanh(t) = \frac{\exp(t) - \exp(-t)}{\exp(t) + \exp(-t)}. \quad (3.51)$$

4097 A curve for a logistic function is shown in Fig. 3.13. The curve has a property that
 4098 $f(-\infty) = 0$ and $f(\infty) = 1$. This property makes the logistic activation function
 4099 useful for categorical assignment: False for 0 and True for 1.

4100 Figure 3.13 can be plotted by the following computer code.

```

4101 #R plot Fig. 9.13: Curve of a logistic function
4102 y = seq(-2, 4, len = 101)
4103 k = 3.2
4104 y0 = 1
4105 setEPS() #Automatically saves the .eps file
4106 postscript("fig0913.eps", height=5, width=7)
4107 par(mar = c(4.2, 4.2, 2, 0.5))
4108 plot(y, 1/(1 + exp(-k*(y - y0))),
4109       type = 'l', col = 'blue', lwd = 2,
4110       xlab = 'y', ylab = 'f(y)',
4111       main = 'Logistic function for k = 3.2, y0 = 1.0',
4112       cex.lab = 1.3, cex.axis = 1.3)
4113 dev.off()
```

```

#Python plot Fig. 9.13: Curve of a logistic function
# define variables
y = np.linspace(-2,4,101)
k = 3.2
y0 = 1

# plot figure
plt.plot(y, 1/(1+ np.exp(-k*(y-y0))), 'b', linewidth = 2.5)
# add labels
plt.title('Logistic function for k=3.2, z0=1.0',
           pad = 10)
plt.xlabel("z", labelpad = 10)
plt.ylabel("g(z)", labelpad = 10)
plt.show()

```

4114

When the growth rate $k = 1$ and midpoint $y_0 = 0$, the logistic function is often referred to as a sigmoid function. Some Python or R code uses sigmoid to represent the activation function.

The idea of data-based decision process is that you integrate all the data, develop a model through data fitting with a specified optimization, and apply the model to new data to generate predictions. The neural network integrates data by assigning each datum a weight as Eq. (3.48) and assigning each weighted sum a bias. The weighted data z_j and the decision data, 0 or 1, from the training dataset form a series of data pairs (z, d) . These data pairs are used for the logistic data fitting many times to train a neural network model. Each time, the NN algorithm will try to optimize its data fitting, such as minimizing mean square errors (MSE) of the fitting. This optimization updates the weights and biases. This optimization process is called back-propagation in an NN algorithm. This makes the NN learning an optimization process. When the fitting errors do not change much, the optimization process has converged, and a logistic function model has been trained. You can now plot the trained neural network, denoted by `nn` in the above code, and produce Fig. 3.12.

The new data are aggregated together as a data frame denoted by `test` in the above code. You apply the trained `nn` model to `test` and obtain the NN prediction result as probabilities, a higher probability suggesting a positive decision, i.e., recruited. You may use 0.5 as your hiring probability threshold and convert the probability result into a decision result, indicated by 1(recruit) or 0 (reject). In this example, Candidates A and B are hired, and Candidate C is rejected, suggested by this NN model.

When running the computer code for this example of hiring data, you may have noticed that each run has a different prediction result. This makes the prediction unreliable. To improve the situation, we need more training data and more optimizations by using several layers of neurons. The approach of multiple hidden layers is a type of deep learning. It is intuitive why more training data are needed to train a more reliable model, since few data cannot cover all possibilities. In the next sub-section, we will present an example of the NN deep learning using more

4145 than six groups of training data. Specifically, we will use 75 groups of the Fisher
4146 iris flower data.

4147 3.4.1.3 Dissect NN results and errors

4148 Logistic regression fits weighted data to a logistic function. Logistic regression can
4149 be a standalone statistics course or a chapter in a statistics book. Interested readers
4150 are referred to those books and courses.

4151 In addition, there are other types of activation functions, such as linear, rectified
4152 linear, leaky rectified linear, exponential linear, and SoftPlus. You can find these
4153 functions from the specialized ML books or online.

4154 Usually a software package of neural network can output the weights and bias,
4155 such as using the following R command

```
4156 nn$weights[[1]][[1]]  
4157
```

4158 in the R NN package **neuralnet**. R can also output the initial assignment of weights
4159 and bias

```
4160 nn$weights[[1]][[2]]  
4161
```

4162 The initial weights and biases are randomly assigned, which makes the NN result
4163 different for different runs when having insufficient amount of data, or insufficient
4164 number of hidden layers.

4165 An NN computer package can also output many other modeling results, such as
4166 an R command

```
4167 nn$result.matrix  
4168
```

4169 These output data can be used to analyze the quality of your trained model. So, it
4170 is unfair to say that NN is a blackbox. Yet, analyzing a trained model for a complex
4171 NN is very difficult and requires some nontrivial mathematical preparation.

4172 Most NN users may not have the mathematical ability to analyze the trained
4173 model based on mathematical theories. Instead, they apply the model to the new
4174 data and see if the NN prediction makes sense based on their experience or common
4175 sense. If not, train the model again with different parameters, or even feed more
4176 training data.

4177 3.4.2 An NN prediction of iris species

4179 We wish to use the Fisher iris data to show an example of NN with a sufficiently
4180 large size of training data and a highly reliable result. The Fisher iris data is a

4181 150 × 5 matrix with three iris species: 50 setosa, 50 virginica, and 50 versicolor, as
 4182 described in the random forest section (Fisher 1936). We wish to use a percentage of
 4183 the data as the training data and the rest as the test data. The following computer
 4184 code is for 50% of the data to be the training data, i.e., $p = 0.5$ in the code. We
 4185 choose 10 neurons, and 10 hidden layers. The NN prediction result shows that NN
 4186 correctly predicted all the 27 setosa irises, got 19 correct among the 20 versicolor
 4187 irises but got one versicolor incorrectly as virginica, and got 26 correct among the
 4188 28 virginica irises but got 2 incorrectly as versicolor. Since the 75 groups of training
 4189 data are plenty, compared to the 6 groups of training data for the hiring decision,
 4190 the trained NN model for the iris flowers does not vary much for different runs.
 4191 The prediction results are quite reliable. Using more neurons and hidden layers may
 4192 also help achieve reliable results. Although the results from different runs of the NN
 4193 code do not change much, small differences still exist. For example, you may get a
 4194 result of predicting 27 virginica.

4195 When you change $p = 0.8$, you will have 120 rows of training data and 30 rows of
 4196 test data. This is the same as what we used in the random forest example section.
 4197 This increase of training data size helps further stabilize the prediction results, i.e.,
 4198 the results from different runs have little differences. You can make some runs of
 4199 the code on your own computer, and see how many incorrect predictions are there.
 4200 Our example runs shown below show only one or two incorrect predictions among
 4201 the 30 rows of test data. For example, in one run, the 9 setosa and 12 versicolor
 4202 irises were correctly predicted, and the 9 virginica irises had 8 correct predictions
 4203 and one wrong. The level of the NN prediction accuracy seems similar to that of
 4204 the random forest prediction for this particular dataset.

```
4205 #R NN code for the Fisher iris flower data
4206 #Ref: https://rpubs.com/vitorhs/iris
4207 data(iris) #150-by-5 iris data
4208 #attach True or False columns to iris data
4209 iris$setosa = iris$Species == "setosa"
4210 iris$virginica = iris$Species == "virginica"
4211 iris$versicolor = iris$Species == "versicolor"
4212 p = 0.5 # assign 50% of data for training
4213 train.idx = sample(x = nrow(iris), size = p*nrow(iris))
4214 train = iris[train.idx,] #determine the training data
4215 test = iris[-train.idx,] #determine the test data
4216 dim(train) #check the train data dimension
4217 #[1] 75 8
4218
4219 #training a neural network
4220 library(neuralnet)
4221 #use the length, width, True and False data for training
4222 iris.nn = neuralnet(setosa + versicolor + virginica ~
4223   Sepal.Length + Sepal.Width +
4224   Petal.Length + Petal.Width,
4225   data = train, hidden=c(10, 10),
4226   rep = 5, err.fct = "ce",
4227   linear.output = F, lifesign = "minimal",
4228   stepmax = 1000000, threshold = 0.001)
4229
```

```
4230 plot(iris.nn, rep="best") #plot the neural network
4231
4232 #Prediction for the rest data
4233 prediction = neuralnet::compute(iris.nn, test[,1:4])
4234 #prediction$net.result is 75-by-3 matrix
4235 prediction$net.result[1:2,]
4236 #[,1] [,2] [,3]
4237 #2 1 3.531638e-09 6.745692e-137
4238 #4 1 3.527872e-09 6.756521e-137
4239
4240 #find which column is for the max of each row
4241 pred.idx <- apply(prediction$net.result, 1, which.max)
4242 pred.idx
4243 #2 4 6 8 9 10
4244 #1 1 1 1 1 1
4245
4246 #Assign 1 for setosa, 2 for versicolor, 3 for virginica
4247 predicted <- c('setosa', 'versicolor', 'virginica')[pred.idx]
4248 predicted[1:6] #The prediction result
4249 #[1] "setosa" "setosa" "setosa" "setosa" "setosa"
4250
4251 #Create confusion matrix: table(prediction,observation)
4252 table(predicted, test$Species)
4253 #predicted      setosa versicolor virginica
4254 #setosa          27       0       0
4255 #versicolor      0       19       2
4256 #virginica      0        1       26
```


4258

3.5 Chapter summary

4259

4260 This chapter has described four methods of machine learning: K-means, support
4261 vector machine (SVM), random forest (RF), and neural network (NN). In the de-
4262 scription, we have used the following datasets: the daily weather data at the Miami
4263 International Airport, the United States, the daily air quality data of New York
4264 City, and R. A. Fisher's iris flower data of species. We have provided both R and
4265 Python codes for these algorithms and their applications examples. The following
4266 provides a brief summary of our descriptions about the four methods.

- 4267 (i) K-means clustering: Simply speaking, this clustering method can fairly divide
4268 N identical candies on a table for K kids according to the candy locations.
4269 It is normally used as an unsupervised learning method. The K-means algo-
4270 rithm minimizes the total within cluster sum of squares (tWCSS) through
4271 iterations. As an example, we applied the K-means clustering method to
4272 the data of the daily minimum surface air temperature T_{min} [in $^{\circ}\text{C}$] and
4273 the direction of the fastest 2-minute wind in a day $WDF2$ [in degrees] of the
4274 Miami International Airport. Our tests suggested that the data have two
4275 clusters, one corresponding to the prevailing wind from east, and another
4276 from west.
- 4277 (ii) Support vector machine for the maximum separation of different sets: The
4278 SVM algorithm is built on the principle of the maximum differences among
4279 the labelled groups of data. The maximum separation of different groups is
4280 measured by the distance between positive and negative hyperplanes. SVM
4281 can also predict the labels of the new unlabeled data. SVM is usually used
4282 as a tool of supervised learning.
- 4283 (iii) Random forest of decision trees: An RF model is a set of decision trees which
4284 form a “forest of decisions.” It is usually used as a tool of supervised learn-
4285 ing. It has two steps: training and prediction. RF can make both classifi-
4286 cation and regression. A benchmark RF example is used in our RF classi-
4287 fication description: The separation and prediction of the R.A. Fisher iris
4288 species dataset. We have also included an example on RF regression using
4289 the daily air quality data of New York City.
- 4290 (iv) Neural network of multiple data fitting: An NN model is a series of data
4291 fitting for a given activation function. NN can also make classifications and
4292 numerical predictions as a regression. The logistic function is a popular
4293 activation function, a smooth transition between category 0 to category 1.
4294 Because of randomness is involved in the NN algorithm, different runs of R
4295 or Python code may yield different results. We have included an example
4296 of using the NN model to help make a hiring decision using the existing
4297 data. To compare with the RF model, we have also applied the NN model
4298 to the R.A. Fisher flower dataset of iris species.

References and Further Readings

- 4300 [1] Boehmke, B., and B. Greenwell, 2019: *Hands-on Machine Learning with R*.
 4301 Chapman and Hall/CRC, Boca Raton, USA.

This machine learning book has many R code and examples. The book website is <https://bradleyboehmke.github.io/HOML/>.

- 4302 [2] Chambers, J. M., W.S. Cleveland, B. Kleiner, and P. A. Tukey, 1983: *Graphical Methods for Data Analysis*. Wadsworth International Group, Belmont, USA.

This book contains the daily ozone and weather data from 1 May 1973 to 30 September 1973 of New York. This dataset serves as a benchmark dataset in random forest algorithm. John M. Chambers is a Canadian statistician, who developed the S programming language, and is a core member of the R programming language project.

- 4305 [3] Fisher, R.A., 1936: The use of multiple measurements in taxonomic problems.
 4307 *Annals of Eugenics*. 7(2), 179-188.

This paper contained the iris data frequently used in the teaching of machine learning. The current name of the journal *Annals of Eugenics* is *Annals of Human Genetics*. The author, Ronald Aylmer Fisher (1890-962), was a British mathematician, statistician, biologist, and geneticist. Modern description of this dataset and its applications in statistics and machine learning can be found from numerous websites, such as

https://www.angela1c.com/projects/iris_project/the-iris-dataset/

- 4308 [4] Géron, A., 2017: *Hands-on Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O'Reilly, Boston, USA.

This machine learning book has many Python code and examples.

- 4312 [5] McCulloch, W. and W. Pitts, 1943: A logical calculus of ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics* 5, 115 -133.

4315

This seminal work is considered the first paper that created the modern NN theory. The paper is highly mathematical, contains ten theorems, and includes a figure of neurons.

4316

- [6] Hastie, T., R. Tibshirani, and J.H. Friedman, 2017: *The Elements of Statistical Learning: Data mining, Inference, and Prediction*. 2nd Ed., Springer, New York, USA.

4317

4318

4319

This is a very famous machine learning book, which has in depth descriptions of the commonly used ML algorithms.

4320

4321

Exercises

4322

- 3.1** Use the K-means method to make a cluster analysis for the following five points in the 2D space: $P_1(1, 1)$, $P_2(2, 2)$, $P_3(2, 3)$, $P_4(3, 4)$, $P_5(4, 4)$. Assume $K = 2$. Plot a figure similar to Fig. 9.1. What is the final tWCSS equal to?

4323

- 3.2** Given the following three points $P_1(1, 1)$, $P_2(2, 2)$, $P_3(2, 3)$ and given $K = 2$, make a K-means cluster analysis following the method presented in Sub-Section 9.1.1. Plot your K-means clustering result.

4324

- 3.3** Use the K-means method to make a cluster analysis for the daily TMIN and WDF2 data of the Miami International Airport in 2015, following the method presented in Sub-sections 9.1.3. Here, TMIN is the daily minimum temperature, and WDF2 denotes the direction [in degrees] of the fastest 2-minute wind in a day. You can obtain the data from the NOAA Climate Data Online website, or from the file

4325

`MiamiIntlAirport2001_2020.csv`

4326

in the `data.zip` file downloadable from the book website

4327

www.climatestatistics.org

4328

- 3.4** Use the K-means method to make a cluster analysis for the daily TMAX and WDF2 data of the Miami International Airport in 2015.

4329

- 3.5** Use the K-means method to make a cluster analysis for the daily TMIN and Tmax data of the Miami International Airport in 2015.

4330

- 3.6** Use the K-means method to make a cluster analysis for the daily TMIN and PRCP data of the Miami International Airport in 2015. Here, PRCP denotes the daily total precipitation.

4331

4332

4333

4334

4335

4336

4337

- 4345 **3.7** Identify two climate parameters of your interest, find the corresponding data
 4346 online, and make a K-means cluster analysis similar to the previous problem
 4347 for your data.
- 4348 **3.8** Following the method presented in Sub-section 9.2.1., make an SVM analysis
 4349 for the following five points in a 2D space: $P_1(1, 1)$, $P_2(2, 2)$, $P_3(2, 3)$, $P_4(3, 4)$, $P_5(4, 4)$.
 4350 The first three points are labeled 1 and the last two are labeled 2. What are
 4351 w , b and D_m ? What points are the support vectors?
- 4352 **3.9** Two new points are introduced to the previous problem: $Q_1(1.5, 1)$ and $Q_2(3, 3)$.
 4353 Use the SVM trained in the previous problem to find out which point belongs
 4354 to what category.
- 4355 **3.10** From the Internet, download the historical daily data of minimum tempera-
 4356 ture (TMIN) and average wind speed (AWND) for a month and a location of
 4357 your interest. Label your data as 1 if the daily precipitation is greater than
 4358 0.5 millimeter, and 0 otherwise. Make an SVM analysis for your labeled data.
- 4359 **3.11** SVM forecast: From the Internet, download the historical daily data of min-
 4360 imum temperature and sea level pressure for a month and a location of your
 4361 interest. Label your data as 1 if the total precipitation of the next day is
 4362 greater than 0.5 millimeter, and 0 otherwise. Make an SVM analysis for your
 4363 labeled data. Given the daily data of minimum temperature and sea level pres-
 4364 sure of a day in the same month but a different year, use your trained SVM
 4365 model to forecast whether the next day was rainy or not, i.e., to determine
 4366 whether the next day is 1 or 0.
- 4367 **3.12** In order to improve the accuracy of your prediction, can you use the data of
 4368 multiple years to train your SVM model? Perform numerical experiments and
 4369 show your results.
- 4370 **3.13** The first 50 in the 150 rows of the R.A. Fisher iris data are for setosa, 51-100
 4371 are for versicolor, and 101-150 are for virginica, use the data 1-40, 51-90, and
 4372 101-140 to train an RF model, following the method in Sub-Section 9.3.1.
 4373 Use the RF model to predict the species of the remaining data of lengths
 4374 and widths of petal and sepal. You can download the R.A. Fisher dataset
 4375 `iris.csv` from the book website or use the data already built in R or Python
 4376 software packages.
- 4377 **3.14** For the 150 rows R.A. Fisher iris data, use only 20% of the data from each
 4378 species to train your RF model. Select another 10% of the riris data of lengths
 4379 and widthes as the new data for prediction. Then use the RF model to predict
 4380 the species of the new data. Discuss the errors of your RF model and your
 4381 prediction.
- 4382 **3.15** Plot a decision tree like Fig. 3.11 for the previous exercise problem.
- 4383 **3.16** For the same R.A. Fisher dataset, design your own RF training and predic-
 4384 tion. Discuss the RF prediction accuracy for this problem.
- 4385 **3.17** RF forecast: From the Internet, download the historical daily data of mini-
 4386 mum temperature and sea level pressure for a month and a location of your
 4387 interest. Label your data as 1 if the total precipitation of the next day is
 4388 greater than 0.5 millimeter, and 0 otherwise. Make an RF analysis for your

- 4389 labeled data. Given the daily data of minimum temperature and sea level
4390 pressure of a day in the same month but a different year, use your trained
4391 RF model to forecast whether the next day was rainy or not. Is your forecast
4392 accurate? How can you improve your forecast?
- 4393 **3.18** Find a climate data time series with missing data and fill in the missing data
4394 using the RF regression method described in Sub-Section 9.3.2.
- 4395 **3.19** The first 50 in the 150 rows of the R.A. Fisher iris data are for setosa, 51-100
4396 are for versicolor, and 101-150 are for virginica, use the data 1-40, 51-90, and
4397 101-140 to train an NN model, following the method in Sub-Section 9.4.2.
4398 Use the NN model to predict the species of the remaining data of lengths and
4399 widths of petal and sepal.
- 4400 **3.20** For the 150 rows R.A. Fisher iris data, use only 20% of the data from each
4401 species to train your NN model. Select another 10% of the riris data of lengths
4402 and widthes as the new data for prediction. Then use the NN model to predict
4403 the species of the new data. Discuss the errors of your NN model and your
4404 prediction.
- 4405 **3.21** NN forecast: From the Internet, download the historical daily data of mini-
4406 mum temperature and sea level pressure for a month and a location of your
4407 interest. Label your data as 1 if the total precipitation of the next day is
4408 greater than 0.5 millimeter, and 0 otherwise. Make an NN analysis for your
4409 labeled data. Given the daily data of minimum temperature and sea level
4410 pressure of a day in the same month but a different year, use your trained
4411 NN model to forecast whether the next day was rainy or not. Is your forecast
4412 accurate? Compare your NN forecast with your RF and SVM forecasts.
- 4413 **3.22** Design a machine learning project for yourself or others. What are your train-
4414 ing data? What are your test data? What is your training model? What is
4415 your training model error? How would you assess your prediction error?

The word “regression” means “a return to a previous and less advanced or worse form, state, condition, or way of behaving,” according to the Cambridge dictionary. The first part “regress” of the word originates from the Latin “regressus,” past participle of regredi (“to go back”), from re- (“back”) + gradi (“to go”). Thus, “regress” means “return, to go back” and is in contrast to the commonly used word “progress.” The *regression* in statistical data analysis refers to a process of returning from the irregular and complex data to a simpler and less perfect state, which is called a model and can be expressed as a curve, a surface, or a function. The function or curve, less complex or less advanced than the irregular data pattern, describes a way of behaving or a relationship. This chapter covers linear models in both uni- and multivariate regressions, least-square estimations of parameters, confidence intervals and inference of the parameters, and fittings of polynomials and other nonlinear curves. By running diagnostic studies on residuals we explain the assumptions of a linear regression model: linearity, homogeneity, independence, and normality. As usual, we use examples of real climate data and provide both R and Python codes.

4.1 Simple linear regression

The simplest model of regression is a straight line, i.e., a linear model. This involves only two variables x and y . Therefore, a *simple linear regression* means going from the data pairs $(x_i, y_i)(i = 1, 2, \dots, n)$ on the xy -plane back to a simple straight line model

$$y = a + bx. \quad (4.1)$$

The data will determine the values of a and b by the criterion of the best fit. The model can be used for prediction: predict y when a value of x is given.

4.1.1 Temperature lapse rate and an approximately linear model

It is known that as an approximation, in a mountain terrain temperature decreases linearly according to elevation. The decrease rate is called the temperature *lapse rate* in meteorology. The linear relationship and the lapse rate are shown in Fig. 4.1 for Colorado, which is a Rocky Mountains state of the United States and has

4449 a large elevation range from 1010 to 4401 meters, according to Colorado Tourism
 4450 www.colorado.com. The dots in the figure are the scatter plot of the data of eleva-
 4451 tion and 1981-2010 average July surface air temperature Tmean of the 24 Colorado
 4452 stations in the United States Historical Climatological Network (USHCN) (Menne
 4453 et al. 2009). The straight line of the figure corresponds to the linear model

$$y = 33.48 - 0.0070x. \quad (4.2)$$

4454 This equation is a result from regression of data, referred to as observations in
 4455 climate science, or sample data in statistics.

4456 Thus, Figure 4.1 shows an exhibit of a scattered and complex data pattern re-
 4457 gressed to a simpler form, in this case, a straight line model. This is called the simple
 4458 linear regression process. The regression result has climate science interpretations,
 4459 such as the lapse rate of temperature.

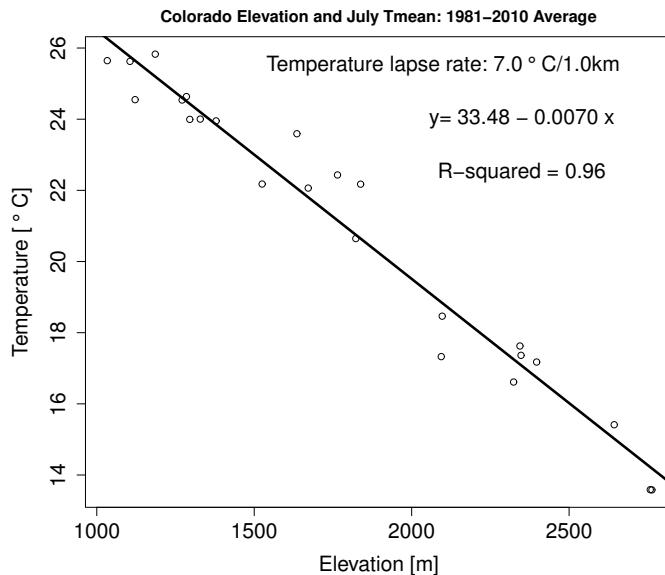


Fig. 4.1 The dots are the scatter plot of the data of elevation and 1981-2010 average July surface air temperature Tmean taken from the 24 USHCN stations. The thick straight line is the linear regression model computed from the data.

4460 The horizontal axis of Fig. 4.1 is for elevation. The elevation data in units of
 4461 meters for the 24 stations are as follows:

4462

4463 1671.5, 1635.6, 2097.0, 1295.4, 1822.7, 2396.9, 2763.0, 1284.7,
 4464 1525.2, 1328.6, 1378.9, 2323.8, 2757.8, 1033.3, 1105.5, 1185.7,
 4465 2343.9, 1764.5, 1271.0, 2347.3, 2094.0, 2643.2, 1837.9, 1121.7

4466 The vertical axis is for temperature, which is the 30-year average of the July daily

4467 mean temperature (Tmean) from 1981-2010, computed from the USHCN monthly
 4468 data, which have been adjusted for the time of observation bias (TOB). Some
 4469 stations had missing data denoted by -9999. When computing the 30-year average,
 4470 the entries of -9999 were omitted. Thus, some averages were computed from fewer
 4471 than 30 years. The resulting average temperature data in the unit of °C for the 24
 4472 stations are as follows:

4473 22.064, 23.591, 18.464, 23.995, 20.645, 17.175, 13.582, 24.635,
 4474 22.178, 24.002, 23.952, 16.613, 13.588, 25.645, 25.625, 25.828,
 4475 17.626, 22.433, 24.539, 17.364, 17.327, 15.413, 22.174, 24.549

4476 With the above data, Figure 4.1 can be generated by the following computer code

```
4477 #R plot Fig. 4.1: Colorado temperature lapse rate
4478 x= c(
4479 1671.5, 1635.6, 2097.0, 1295.4, 1822.7, 2396.9, 2763.0, 1284.7,
4480 1525.2, 1328.6, 1378.9, 2323.8, 2757.8, 1033.3, 1105.5, 1185.7,
4481 2343.9, 1764.5, 1271.0, 2347.3, 2094.0, 2643.2, 1837.9, 1121.7)
4482 y= c(
4483 22.064, 23.591, 18.464, 23.995, 20.645, 17.175, 13.582, 24.635,
4484 22.178, 24.002, 23.952, 16.613, 13.588, 25.645, 25.625, 25.828,
4485 17.626, 22.433, 24.539, 17.364, 17.327, 15.413, 22.174, 24.549)
4486 setEPS() # save the .eps figure
4487 postscript("fig0401.eps", width = 8)
4488 par(mar=c(4.5,4.5,2.5,0.5))
4489 plot(x,y,
4490   xlab="Elevation [m]",
4491   ylab=expression("Temperature [^degree~"C]"),
4492   main="Colorado Elevation and July Tmean: 1981-2010 Average",
4493   cex.lab=1.5, cex.axis=1.5, cex.main =1.2)
4494 reg=lm(y~x)
4495 reg
4496 #(Intercept)           x
4497 # 33.476216      -0.006982    #-7.0 degC/1km
4498 summary(reg)
4499 #R-squared:  0.9631
4500 abline(reg,lwd=3)
4501 text(2100, 25.5,
4502 expression("Temperature_lapse_rate: 7.0"~degree~"C/1.0km"),
4503 cex=1.5)
4504 text(2350, 24, "y= 33.48 - 0.0070 x", cex=1.5)
4505 text(2350, 22.5,"R-squared= 0.96", cex=1.5)
4506 dev.off()
```

```

# Python plot Fig. 4.1: Colorado temperature lapse rate
x = np.array([
    1671.5, 1635.6, 2097.0, 1295.4, 1822.7, 2396.9, 2763.0, 1284.7,
    1525.2, 1328.6, 1378.9, 2323.8, 2757.8, 1033.3, 1105.5, 1185.7,
    2343.9, 1764.5, 1271.0, 2347.3, 2094.0, 2643.2, 1837.9, 1121.7
])
y = np.array([
    22.064, 23.591, 18.464, 23.995, 20.645, 17.175, 13.582, 24.635,
    22.178, 24.002, 23.952, 16.613, 13.588, 25.645, 25.625, 25.828,
    17.626, 22.433, 24.539, 17.364, 17.327, 15.413, 22.174, 24.549
])
# calculate correlation coefficients
corrMatr = np.corrcoef(x,y)
# R-squared
Rsqu = corrMatr[0,1]**2
# trend line
reg1 = np.array(np.polyfit(x, y, 1))
abline = reg1[1] + x*reg1[0]
fig, ax = plt.subplots(figsize=(12,12))
ax.plot(x,y, 'ko');
ax.plot(x,abline, 'k-');
ax.set_title("Colorado\u2022Elevation\u2022vs.\u2022July\u2022Tmean:\u2022\n\u20221981\u2022\u20222010\u2022Average",
             size = 25, fontweight = 'bold', pad = 20)
ax.set_xlabel("Elevation\u2022[$m$]", size = 25, labelpad = 20)
ax.set_ylabel("Temperature\u2022[$\degree C$]", size = 25, labelpad = 20);
ax.tick_params(length=6, width=2, labelsize=20);
ax.set_xticks(np.round(np.linspace(1000, 3000, 5), 2))
ax.set_yticks(np.round(np.linspace(12, 26, 8), 2))
ax.text(1750, 25.5,
        r"Temperature\u2022lapse\u2022rate:\u20227.0\u2022$\degree C/km$",
        color= 'k', size = 20)
ax.text(2250, 24.5, r"$y\u2022=\u202233.48\u2022-\u20220.0070\u2022x$" % Rsqu,
        color= 'k', size = 20)
ax.text(2250, 23.5, r"$R\u2022-squared\u2022=\u2022%.2f$" % Rsqu,
        color= 'k', size = 20)
plt.show()

```

4507

4508 The Colorado July temperature lapse rate (TLR) $7.0 \text{ } ^\circ\text{C}/1.0\text{km}$ is consistent with
 4509 earlier studies using 104 stations on the west slope of Colorado for an elevation range
 4510 of 1,500- 3,600 meters: $6.9 \text{ } ^\circ\text{C}/1.0\text{km}$ (Fall 1997). For the annual mean temperature,
 4511 the LTR is $6.0 \text{ } ^\circ\text{C}/1.0\text{km}$. These two results are comparable to an LTR study for the
 4512 Southern Ecuadorian Andes in the elevation range of 2,610-4,200 meters (Cordova
 4513 et al. 2016). The annual mean temperature TLR is $6.88 \text{ } ^\circ\text{C}/1.0\text{km}$.

4514 TLR may be applied to approximately predict the temperature of a mountain
 4515 region at a given location, in case it is hard to maintain a weather station at some
 4516 locations due to high elevation or complex terrain, while it is relatively easy to
 4517 obtain the elevation data based on the Geographical Information System (GIS) or
 4518 a digital elevation model (DEM) dataset.

4519 **4.1.2 Assumptions and formula derivations of the single variate**
 4520 **linear regression**

4522 This subsection describes the statistical concepts and mathematical theory of the
 4523 above linear regression procedure. We pay special attention to the assumptions
 4524 about the linear regression model.

4525 **4.1.2.1 Linear model and its assumptions**

4526 The linear model for TLR may be written as

$$Y = a + bx + \epsilon. \quad (4.3)$$

4527 Here,

- 4528 (i) x is an independent variable, also called explanatory variable, and is deterministic (not random);
- 4529 (ii) Y is the dependent variable, also called the response variable, and is a random variable with a constant variance $\text{Var}[Y|x] = \sigma^2$ for any x ;
- 4530 (iii) ϵ is the random error term, which is assumed to have zero mean $E[\epsilon] = 0$, and constant variance σ^2 : $\text{Var}(\epsilon|x) = \sigma^2$, and is uncorrelated with each other: $\text{Cor}(\epsilon|x_1, \epsilon|x_2) = 0$ if $x_1 \neq x_2$;
- 4531 (iv) a, b and σ are called parameters and are to be estimated from data (x_i, y_i) ($i = 1, 2, \dots, n$), also called samples, or sample data; and
- 4532 (v) The expected value of Y for a given x is $a + bx$, i.e.,

$$E[Y|x] = a + bx. \quad (4.4)$$

4538 The linear model Eq. (4.3) at the given points x_i is

$$Y_i = a + bx_i + \epsilon_i, \quad i = 1, 2, \dots, n. \quad (4.5)$$

4539 where both ϵ_i and Y_i are random variables with

$$E[Y_i] = a + bx_i, \quad (4.6)$$

$$E[\epsilon_i] = 0, \quad (4.7)$$

$$\text{Var}(Y_i) = E[(Y_i - (a + bx_i))^2] = \sigma^2, \quad (4.8)$$

$$\text{Var}(\epsilon_i) = E[(\epsilon_i - 0)^2] = \sigma^2, \quad (4.9)$$

$$\text{Cor}(\epsilon_i, \epsilon_j) = \delta_{ij}, \quad (4.10)$$

4540 where δ_{ij} is the Kronecker delta

$$\delta_{ij} = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{if } i \neq j. \end{cases} \quad (4.11)$$

4541 The last equation means that the error terms ϵ_i at the observational points x_i are
 4542 uncorrelated with each other.

4543 The observations at the given points x_i are y_i . Both are fixed values, such as tem-
 4544 perature $y_1 = 22.064^\circ\text{C}$ at the given elevation $x_1 = 1671.5$ meters in the Colorado
 4545 July temperature example. Here, y_1 is a sample value for the random variable Y_1 .
 4546 The corresponding error datum is e_1 , which is called a residual and is regarded a
 4547 sample value of ϵ_1 with respect to the linear model.

4548 4.1.2.2 Estimating a and b from data

4549 Given a set of sample data $(x_i, y_i), i = 1, 2, \dots, n$, the parameters a and b in the
 4550 regression model can be estimated as \hat{a}, \hat{b} . Thus, \hat{a}, \hat{b} correspond to a particular
 4551 dataset, i.e., another dataset $(u_i, v_i), i = 1, 2, \dots, m$ will yield a pair of different es-
 4552 timators that may be denoted by $\hat{a}_{uv}, \hat{b}_{uv}$. In contrast, a, b are the general notations
 4553 of constant coefficients for a linear regression model.

4554 The estimated linear model may be written in the following way:

$$y_i = \hat{a} + \hat{b}x_i + e_i, \quad i = 1, 2, \dots, n, \quad (4.12)$$

4555 where

$$e_i = y_i - (\hat{a} + \hat{b}x_i) \quad (4.13)$$

4556 are called residuals of the linear model, and are equal to the observed values y_i
 4557 minus the predicted values

$$\hat{y}_i = \hat{a} + \hat{b}x_i. \quad (4.14)$$

4558 The mean of the above n equations (4.12) yields

$$\bar{y} = \hat{a} + \hat{b}\bar{x}, \quad (4.15)$$

4559 where

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}, \quad (4.16)$$

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{n}. \quad (4.17)$$

4560 Because of the unbiased model assumption $E[\epsilon] = 0$, we require that

$$\bar{\epsilon} = \frac{e_1 + e_2 + \dots + e_n}{n} = 0, \quad (4.18)$$

4561 equivalently,

$$\sum_{i=1}^n e_i = 0. \quad (4.19)$$

4562 Equations (4.15) and (4.18) give an estimate for a from data:

$$\hat{a} = \bar{y} - \hat{b}\bar{x}, \quad (4.20)$$

4563 but \hat{b} is to be estimated by

$$\hat{b} = \frac{\mathbf{y}_a \cdot \mathbf{x}_a}{\mathbf{x}_a \cdot \mathbf{x}_a}. \quad (4.21)$$

⁴⁵⁶⁴ Here,

$$\mathbf{x}_a = \begin{bmatrix} x_{a,1} \\ x_{a,2} \\ \vdots \\ x_{a,n} \end{bmatrix} \quad \mathbf{y}_a = \begin{bmatrix} y_{a,1} \\ y_{a,2} \\ \vdots \\ y_{a,n} \end{bmatrix} \quad \mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}, \quad (4.22)$$

⁴⁵⁶⁵ with

$$x_{a,i} = x_i - \bar{x}, \quad (4.23)$$

$$y_{a,i} = y_i - \bar{y}. \quad (4.24)$$

⁴⁵⁶⁶ In the last two formulas, $x_{a,i}$ and $y_{a,i}$ may be regarded as anomaly data, hence
⁴⁵⁶⁷ with the subscript “a.” Climate scientists often uses the concept of anomaly data,
⁴⁵⁶⁸ departures from a climate normal.

⁴⁵⁶⁹ The \hat{b} estimate formula (4.21) can be derived from the best fit condition which
⁴⁵⁷⁰ minimizes the mean square errors (MSE)

$$\text{MSE} = \frac{|\mathbf{e}|^2}{n} = \frac{\sum_{i=1}^n e_i^2}{n} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}, \quad (4.25)$$

⁴⁵⁷¹ equivalently minimizing the sum of the square errors (SSE)

$$\text{SSE} = |\mathbf{e}|^2 = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = n \times \text{MSE}, \quad (4.26)$$

⁴⁵⁷² where $|\mathbf{e}|$ denotes the Euclidean length of vector \mathbf{e} . Thus, this method is also called
⁴⁵⁷³ the least square estimate. The least square condition is equivalent to the following
⁴⁵⁷⁴ orthogonality condition (See Fig. 4.2)¹

$$\mathbf{e} \cdot \mathbf{x}_a = 0. \quad (4.27)$$

⁴⁵⁷⁵

⁴⁵⁷⁶ Inserting Eq. (4.20) into the linear model with data (4.12) leads to

$$y_i - \bar{y} = \hat{b}(x_i - \bar{x}) + e_i, \quad i = 1, 2, \dots, n. \quad (4.28)$$

⁴⁵⁷⁷ or

$$y_{a,i} = \hat{b}x_{a,i} + e_i, \quad i = 1, 2, \dots, n, \quad (4.29)$$

¹ SSE can be written as

$$\text{SSE} = |\mathbf{e}|^2 = |\mathbf{y}_a - \hat{b}\mathbf{x}_a|^2.$$

When \hat{b} is optimized to minimize SSE, the derivative of the above equation with respect to \hat{b} is zero.

$$\frac{d}{d\hat{b}} \text{SSE} = (\mathbf{y}_a - \hat{b}\mathbf{x}_a) \cdot \mathbf{x}_a = \mathbf{e} \cdot \mathbf{x}_a = 0.$$

This is the orthogonality condition. This condition $\mathbf{e} \cdot \mathbf{x}_a = 0$ can also be derived from geometric point of view. The three vectors $\hat{b}\mathbf{x}_a$, \mathbf{e} and \mathbf{y}_a in $\mathbf{y}_a = \hat{b}\mathbf{x}_a + \mathbf{e}$ form a triangle. The side \mathbf{e} , depending on \hat{b} , is the shortest when $\hat{b}\mathbf{x}_a$ is perpendicular to \mathbf{e} by adjusting parameter \hat{b} .

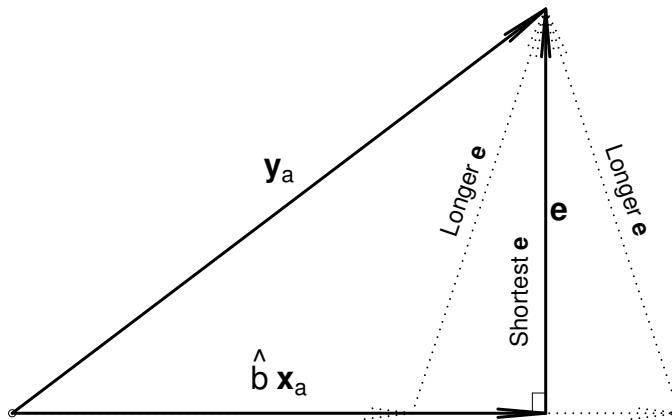


Fig. 4.2 The condition of the least sum of the residual squares SSE is equivalent to the orthogonality condition $\mathbf{e} \cdot \mathbf{x}_a = 0$.

4578 or

$$\mathbf{y}_a = \hat{b}\mathbf{x}_a + \mathbf{e}. \quad (4.30)$$

4579 The dot product of both sides of Eq. (4.29) with \mathbf{x}_a yields

$$\hat{b} = \frac{\mathbf{y}_a \cdot \mathbf{x}_a - \mathbf{e} \cdot \mathbf{x}_a}{\mathbf{x}_a \cdot \mathbf{x}_a}. \quad (4.31)$$

4580 With the orthogonality condition Eq. (4.27) $\mathbf{e} \cdot \mathbf{x}_a = 0$, Eq. (4.31) is reduced to the
4581 least square estimate of b : Eq. (4.21).

4582 In summary, the estimates \hat{a} and \hat{b} are derived by requiring (i) zero mean of the
4583 residuals, and (ii) minimum sum of the residual squares.

4584 As shown in the computer code for Fig. 4.1, the R command for the simple linear
4585 regression is `lm(y ~ x)`. The TLR dataset yields the estimate of the intercept
4586 $\hat{a} = 33.476216$ and slope $\hat{b} = -0.006982$:

```
4587 lm(y ~ x)
4588 #(Intercept)      x
4589 #33.476216     -0.006982
```

4590 The estimated linear model is thus

$$y = 33.476216 - 0.006982x, \quad (4.32)$$

4591 or

$$\text{Temp } [{}^\circ\text{C}] = 33.476216 - 0.006982 \times \text{Elevation } [\text{m}]. \quad (4.33)$$

4592 The corresponding Python code is as follows:

```

#Python linear regression as the first order polynomial fit
reg = np.polyfit(x, y, 1)
print(reg)
#[ -6.98188456e-03  3.34763004e+01]

```

4593

4594 The 24 residuals $e_i, i = 1, 2, \dots, 24$ are

```

4595 reg = lm(y ~ x)
4596 round(reg$residuals, digits = 5)
4597 #0.25792  1.53427 -0.37129 -0.43697 -0.10542  0.43358
4598 #-0.60335  0.12833 -0.64953 -0.19817  0.10302 -0.63880
4599 #-0.63366 -0.61692 -0.13283  0.63012  0.51454  1.27623
4600 #-0.06333  0.27628 -1.52923  0.39122  1.52971 -1.09572

```

4601 The mean of the 24 residuals is zero, demonstrated by an R code below

```

4602 mean(reg$residuals)
4603 #[1] 1.62043e-17

```

4604 The least square condition $\mathbf{e} \cdot \mathbf{x}_a = 0$ is demonstrated by the following R code

```

4605 xa = x - mean(x)
4606 sum(xa*reg$residuals)
4607 #[1] -2.83773e-13

```

4608 The unbiased MSE is

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2. \quad (4.34)$$

4609 Here, subtracting 2 is due to the two constraints of estimating a and b . The unbiased
4610 MSE s^2 the Colorado TLR regression can be computed using the following R code

```

4611 sum((reg$residuals)^2)/(length(y) - 2)
4612 #[1] 0.6096193

```

4613 Thus, one can use the data to estimate \hat{b} first using (4.21), then \hat{a} using (4.20),
4614 and finally $\hat{\sigma}^2$, also denoted by s^2 , using (4.34).

4615 The above computer code is summarized as follows.

```

4616 #R code for the Colorado TLR regression analysis
4617 lm(y ~ x)
4618 #(Intercept)           x
4619 #33.476216   -0.006982
4620
4621 reg = lm(y ~ x)
4622 round(reg$residuals, digits = 5)
4623 mean(reg$residuals)
4624 #[1] 1.62043e-17
4625
4626 xa = x - mean(x)
4627 sum(xa*reg$residuals)
4628 #[1] -2.83773e-13

```

```

# Python code for the Colorado TLR regression analysis
reg = np.polyfit(x, y, 1)
print(reg) #slope and intercept
#[ -6.98188456e-03  3.34763004e+01]

regfit = np.polyval(reg,x)
regresiduals = y - regfit
print(np.round(regresiduals,5))
print(np.mean(regresiduals))
# -1.850371707708594e-17 #should be almost zero

xa = np.array(x) - np.mean(x)
np.sum(xa*regresiduals)
print(np.dot(xa, regresiduals)) #orthogonality
# -1.48929757415317e-11 #should be almost zero

np.sum((regresiduals)**2)/(y.size - 2)
# 0.6096193452251238 #unbiased MSE

```

4629

4630 Once again, we emphasize that two assumptions are used in the derivation of the
 4631 formulas to estimate \hat{a} and \hat{b} :

- 4632 (a) The unbiased model assumption: The mean residual is zero $\bar{e} = 0$, and
 4633 (b) The optimization assumption: The residual vector e is perpendicular to the
 4634 x -anomaly vector $x_a = x - \bar{x}$.

4635 Under these two assumptions, the estimators \hat{a} and \hat{b} are highly sensitive to the
 4636 Y outliers, particularly the outlier data corresponding to the smallest and largest
 4637 x values. One outlier can completely change the \hat{a} and \hat{b} values, in agreement with
 4638 our intuition. This is an end point problem often encountered in data analysis. To
 4639 suppress the sensitivity, many robust regression methods have been developed and
 4640 R packages are available, such as the Robust Regression, by UCLA institute for
 4641 Digital Research & Education:

4642 <https://stats.idre.ucla.edu/r/dae/robust-regression/>

4643 Figure 4.2 can be generated by the following R code.

```

#R plot Fig. 4.2: Geometric derivation of the least squares
par(mar=c(0.0,0.5,0.0,0.5))
plot(0,0, xlim=c(0,5.2), ylim=c(0,2.2),
      axes = FALSE, xlab="", ylab="")
arrows(0,0,4,0, angle=5, code=2, lwd=3, length=0.5)
arrows(4,0,4,2, angle=5, code=2, lwd=3, length=0.5)
arrows(0,0,4,2, angle=5, code=2, lwd=3, length=0.5)
arrows(5,0,4,2, angle=7, code=2, lwd=2, lty=3, length=0.5)
arrows(0,0,5,0, angle=7, code=2, lwd=2, lty=3, length=0.5)
arrows(3,0,4,2, angle=7, code=2, lwd=2, lty=3, length=0.5)
arrows(0,0,3,0, angle=7, code=2, lwd=2, lty=3, length=0.5)
segments(3.9,0, 3.9, 0.1)
segments(3.9, 0.1, 4.0, 0.1)
text(2,0.2, expression(hat(b)^bold(x)[a]), cex=2)
text(2,1.2, expression(bold(y)[a]), cex=2)
text(4.1,1, expression(bold(e)), cex=2)

```

```

4660 text(3.8,0.6, expression(paste("Shortest",bold(e))),  

4661   cex=1.5, srt=90)  

4662 text(3.4,1.1, expression(paste("Longer",bold(e))),  

4663   cex=1.5, srt=71)  

4664 text(4.6,1.1, expression(paste("Longer",bold(e))),  

4665   cex=1.5, srt=-71)

# Python plot Fig. 4.2: Geometry of the least squares
plt.figure(figsize=(10,8))# Define figure size
plt.xlim([0, 5.2]); plt.ylim([0, 2.2])
plt.axis('off')
plt.annotate(' ', ha = 'center', va = 'bottom',
            xytext = (0, 0),xy = (4, 0),
            arrowprops = {'facecolor' : 'black'})
plt.annotate(' ', ha = 'center', va = 'bottom',
            xytext = (4, 0),xy = (4, 2),
            arrowprops = {'facecolor' : 'black'})
plt.annotate(' ', ha = 'center', va = 'bottom',
            xytext = (0, 0),xy = (4, 2),
            arrowprops = {'facecolor' : 'black'})
plt.plot([4, 5], [0, 0], color ='k', linewidth = 12,
         dashes = (3, 1))
plt.plot([3,4], [0, 2], color ='k', linewidth = 3,
         dashes = (3, 1))
plt.plot([5,4], [0, 2], color ='k', linewidth = 3,
         dashes = (4,1))
plt.plot([0, 3], [0, 0], color ='k', linewidth = 12,
         dashes = (4,1))
plt.plot([3.9, 3.9], [0, 0.1], color ='k')
plt.plot([3.9, 4.0], [0.1, 0.1], color ='k')
plt.text(2, 0.1, r'$\hat{b} \mathbf{x}_a$', fontsize = 30)
plt.text(2, 1.2, r'$\mathbf{y}_a$', fontsize = 30)
plt.text(4.08, 0.8, r'$\mathbf{e}$', fontsize = 30)
plt.text(3.8, 0.6, r'Shortest $\mathbf{e}$',
         rotation=90, fontsize = 20)
plt.text(3.3, 1.1, r'Longer $\mathbf{e}$',
         rotation=71, fontsize = 20)
plt.text(4.3, 1.1, r'Longer $\mathbf{e}$',
         rotation= -71, fontsize = 20)
plt.show()

4666

```

4.1.2.3 Relationship between slope and correlation

4667 The slope \hat{b} and correlation r_{xy} are related in the following way

$$\hat{b}|\mathbf{x}_a| = r_{xy}|\mathbf{y}_a|. \quad (4.35)$$

4668 This has two extreme cases: $r_{xy} = 1$ and $r_{xy} = 0$.

4669 Case (a). Perfect correlation $r_{xy} = 1$ when the points on the scatter plot of the x, y data lie exactly on a straight line, which implies perfect prediction

$$\hat{b}|\mathbf{x}_a| = |\mathbf{y}_a|. \quad (4.36)$$

4672 Case (b). Perfect noise with zero correlation $r_{xy} = 0$, which implies

$$\hat{b} = 0, \quad (4.37)$$

4673 i.e., the best prediction is the mean.

4674 Formula (4.35) can be derived as follows. The slope estimate formula (4.21) can
4675 be re-written to provide geometric and science interpretations.

4676 (i) Geometric projection interpretation:

$$\hat{b} = \frac{\mathbf{y}_a \cdot (\mathbf{x}_a / |\mathbf{x}_a|)}{|\mathbf{x}_a|}. \quad (4.38)$$

4677 Since $\mathbf{x}_a / |\mathbf{x}_a|$ is a unit vector in the direction of \mathbf{x}_a , the slope is the projec-
4678 tion of the y anomaly data vector on the x anomaly data vector and then
4679 normalized by the x anomaly data vector.

4680 Or we write the \hat{b} formula in the following way:

$$\hat{b} = \frac{\mathbf{y}_a}{|\mathbf{x}_a|} \cdot \frac{\mathbf{x}_a}{|\mathbf{x}_a|}. \quad (4.39)$$

4681 The first fraction on the right hand side is the y anomaly data vector
4682 normalized by $|\mathbf{x}_a|$. Thus, the slope is the projection of this normalized y
4683 anomaly data vector onto the unit x anomaly data vector.

4684 (ii) Correlation interpretation:

$$\hat{b} = r_{xy} \frac{|\mathbf{y}_a|}{|\mathbf{x}_a|} \quad (4.40)$$

4685 where r_{xy} is the correlation coefficient, or simply called the correlation,
4686 defined by

$$r_{xy} = \frac{\mathbf{y}_a \cdot \mathbf{x}_a}{|\mathbf{x}_a| |\mathbf{y}_a|} \quad (4.41)$$

4687 The slope depends directly on the correlation coefficient between the x and
4688 y anomaly data. The correlation is scaled by the ratio of the length of the
4689 y anomaly data vector to the length of the x anomaly data vector.

4690 Below is the computer code corresponding to the two interpretations for the
4691 Colorado July TLR example.

```
4692 #R code for estimating regression slope b
4693
4694 #Method 1: Using vector projection
4695 xa = x - mean(x) #Anomaly the x data vector
4696 nxa = sqrt(sum(xa^2)) #Norm of the anomaly data vector
4697 ya = y - mean(y)
4698 nya=sqrt(sum(ya^2))
4699 sum(ya*(xa/nxa))/nxa #Compute b
4700 #[1] -0.006981885 #This is an estimate for b
4701
4702 #Method 2: Using correlation
4703 corxy=cor(xa, ya) #Compute the correlation between xa and ya
4704 corxy
4705 #[1] -0.9813858 #Very high correlation
```

```
4706 corxy*nya/nxa #Compute b
4707 #[1] -0.006981885 #This is an estimate for b
```

```
#Python code for estimating regression slope b

#Method 1: Using vector projection
xa = x - np.mean(x) #Anomaly the x data vector
nxa = np.sqrt(np.sum(xa**2)) #Norm of the anomaly vector
ya = y - np.mean(y)
nya = np.sqrt(sum(ya**2))
print(np.sum(ya*(xa/nxa))/nxa) #Compute b
#[1] -0.006981885 #This is an estimate for b

#Method 2: Using correlation
from scipy.stats import pearsonr
corxy, _ = pearsonr(xa, ya) #Correlation between xa and ya
print(corxy)
#[1] -0.9813858 #Very high correlation
print(corxy*nya/nxa) #Compute b
#[1] -0.006981885 #This is an estimate for b
```

4708

4.1.2.4 Percentage of variance explained: $R^2 \times 100\%$

4709 The unbiased variance of the model-predicted data $\hat{y}_i(i = 1, 2, \dots, n)$ is defined as

$$\text{MV} = \frac{\sum_{i=1}^n [(\hat{y}_i - \bar{\hat{y}})]^2}{n - 1}, \quad (4.42)$$

4710 where $\bar{\hat{y}}$ is the mean of $\hat{y}_i(i = 1, 2, \dots, n)$.

4711 The unbiased variance of the Y data

$$\text{YV} = \frac{\sum_{i=1}^n [(y_i - \bar{y})]^2}{n - 1}, \quad (4.43)$$

4712 where \bar{y} is the mean of the original station temperature data $y_i(i = 1, 2, \dots, n)$.

4713 We wish to measure how good the model is in terms of variance, and thus define
4714 the ratio of MV to YV, which is named R^2

$$R^2 = \frac{\text{MV}}{\text{YV}}, \quad (4.44)$$

4715 The value $R^2 \times 100\%$ indicates the percentage of variance explained by the linear
4716 model. A larger R^2 value suggests a better model.

4718 The variance of the model-predicted data can be re-written as

$$\begin{aligned}
 MV &= \frac{\sum_{i=1}^n [(\hat{a} + \hat{b}x_i) - (\hat{a} + \hat{b}\bar{x})]^2}{n-1} \\
 &= \hat{b}^2 \frac{\sum_{i=1}^n [x_i - \bar{x}]^2}{n-1} \\
 &= \left[r_{xy} \frac{|\mathbf{y}_a|}{|\mathbf{x}_a|} \right]^2 \frac{|\mathbf{x}_a|^2}{n-1} \\
 &= r_{xy}^2 \frac{|\mathbf{y}_a|^2}{n-1} \\
 &= r_{xy}^2 \times YV.
 \end{aligned} \tag{4.45}$$

4719 Therefore,

$$R^2 = r_{xy}^2. \tag{4.46}$$

4720 The R^2 value is equal to the square of the correlation coefficient computed from
4721 the data $(x_i, y_i) (i = 1, 2, \dots, n)$ with $0 \leq R^2 \leq 1$.

4722 For the Colorado TLR data used earlier, we have $r_{xy} = -0.9813858$, and $R^2 =$
4723 0.9631181 .

4724 The variances MV, YV, and R^2 can be computed by the following computer code
4725 in multiple ways.

```

4726 #R code for computing MV
4727 var(reg$fitted.values)
4728 #[1] 15.22721
4729 #Or another way
4730 yhat = reg$fitted.values
4731 var(yhat)
4732 #[1] 15.22721
4733 #Or still another way
4734 n = 24
4735 sum((yhat - mean(yhat))^2)/(n-1)
4736 #[1] 15.22721
4737
4738 #R code for computing YV
4739 sum((y - mean(y))^2)/(n-1)
4740 # [1] 15.81033
4741 #Or another way
4742 var(y)
4743 #[1] 15.81033
4744
4745 #R code for computing R-squared value
4746 var(reg$fitted.values)/var(y)
4747 #[1] 0.9631181 #This is the R-squared value
4748
4749 cor(x,y)
4750 #[1] -0.9813858
4751 (cor(x,y))^2
4752 #[1] 0.9631181 #This is the R-squared value

```

```

#Python code for Compute MV, YV, and R^2
from statistics import variance
reg = np.polyfit(x, y, 1)
yhat = np.polyval(reg,x)
#computing MV
print('MV=,', variance(yhat))
#MV = 15.227212262175959
#Or another way
n = 24
print('MV=,', np.sum((yhat - np.mean(yhat))**2)/(n-1))
#MV = 15.227212262175959

#computing YV
print('YV=,', np.sum((y - np.mean(y))**2)/(n-1))
#[1] 15.81033
#Or another way
print('YV=,', variance(y))
#YV = 15.81032641847826

#computing R-squared value
print('R-squared=,', variance(yhat)/variance(y))
#R-squared = 0.9631181456430407

#computing correlation and R-squared
from scipy.stats import pearsonr
corxy, _ = pearsonr(x, y)
print('Correlation r_xy=,', corxy)
#Correlation r_xy = -0.9813858291431772
print('R-squared=,', corxy**2)
#R-squared = 0.9631181456430413

```

4753

4754

4.1.2.5 The regression estimates by the least square or the maximum likelihood

4755

4756 A more complicated and commonly used derivation is based on the minimization
 4757 of the sum of squared errors (SSE)

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n [y_i - (\hat{a} + \hat{b}x_i)]^2. \quad (4.47)$$

4758

This approach can be found in most statistics books and from the Internet. The
 4759 terminology of “least square” regression comes from this minimization principle.

4760

This minimization is with respect to \hat{a} and \hat{b} . The minimization condition leads
 4761 to two linear equations that determine \hat{a} and \hat{b} . The “least square” minimization
 4762 condition for \hat{b} is equivalent to the orthogonality condition (4.27). Geometrically,
 4763 the minimum distance of a point to a line is defined as the length of the line
 4764 segment that is orthogonal to the line and connects the point to the line, that is, it

4765 is the minimum distance between the point and the line (See Fig. 4.2). Thus, the
 4766 orthogonality condition and minimization condition are equivalent.

4767 Another commonly used method to estimate the parameters is to maximize a
 4768 likelihood function defined as

$$L(a, b, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp(-(y_i - (a + bx_i))^2 / (2\sigma^2)). \quad (4.48)$$

4769 The solution for the maximum $L(a, b, \sigma)$ yields the estimate of \hat{a} , \hat{b} , and $\hat{\sigma}^2$. The
 4770 result is exactly the same as the ones obtained by the method of least squares or
 4771 perpendicular projection.

4772 However, the maximum likelihood approach shown here explicitly assumes the
 4773 normal distribution of the error term ϵ_i and the response variable Y_i . For the least
 4774 squares approach, the assumption of normal distribution is needed only when for
 4775 the inference of the parameters \hat{a} , \hat{b} , and $\hat{\sigma}^2$, but is not needed to estimate them.

4776 4.1.3 Statistics of slope and intercept: Distributions, confidence 4777 intervals, and inference

4779 4.1.3.1 Mean and variance of the slope

4780 The normal distribution of a and b is now assumed in this section. This assumption
 4781 was not needed earlier to estimate \hat{a} and \hat{b} and to compute R-squared.

4782 Different datasets will yield different estimates of \hat{a} , \hat{b} , and $\hat{\sigma}^2$, which form a dis-
 4783 tribution corresponding to the random datasets. Thus, we may regard \hat{a} , \hat{b} , and $\hat{\sigma}^2$
 4784 as random variables. Their expected values are a , b and σ^2 if the random datasets
 4785 satisfy the linear model assumptions (linearity, constant variance, independent er-
 4786 rors, and normal distribution). However, in practical applications, we just have one
 4787 dataset, estimate \hat{a} , \hat{b} , and $\hat{\sigma}^2$ once, and then interpret the results.

4788 Following Eq. (4.21), instead of using data y_i to estimate the slope, we use the
 4789 corresponding random variables Y_i to define the slope B as a random variable:

$$\begin{aligned} B &= \frac{\mathbf{Y}_a \cdot \mathbf{x}_a}{\mathbf{x}_a \cdot \mathbf{x}_a} \\ &= \frac{\sum_{i=1}^n (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) Y_i, \end{aligned} \quad (4.49)$$

4790 because

$$\sum_{i=1}^n (x_i - \bar{x}) = 0. \quad (4.50)$$

4791 The above shows that B is a linear combination of n terms Y_i . If $Y_i \sim N(a + bx, \sigma^2)$,
 4792 then B is also normally distributed when n is large, because of the Central Limit
 4793 Theorem.

⁴⁷⁹⁴ The expected value of B is b , as shown below:

$$\begin{aligned}
 E[B] &= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) E[Y_i] \\
 &= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) (a + bx_i) \\
 &= b \times \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) x_i \\
 &= b \times \left(\frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \\
 &= b.
 \end{aligned} \tag{4.51}$$

⁴⁷⁹⁵ The variance of B is

$$\text{Var}(B) = \frac{\sigma^2}{\|\mathbf{x}\|_a^2}. \tag{4.52}$$

⁴⁷⁹⁶ This can be shown as follows:

$$\begin{aligned}
 \text{Var}[B] &= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^2 \text{Var}(Y_i) \\
 &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
 &= \frac{\sigma^2}{\|\mathbf{x}\|_a^2}.
 \end{aligned} \tag{4.53}$$

⁴⁷⁹⁷ Statistics literature often denotes

$$S_{xx} = \|\mathbf{x}\|_a^2. \tag{4.54}$$

⁴⁷⁹⁸ Hence,

$$\text{Var}(B) = \frac{\sigma^2}{S_{xx}}. \tag{4.55}$$

⁴⁷⁹⁹ When n goes to infinity, S_{xx} also goes to infinity. Thus,

$$\lim_{n \rightarrow \infty} \text{Var}(B) = 0. \tag{4.56}$$

⁴⁸⁰⁰ This means that more data points yield a better estimate for the slope, supporting
⁴⁸⁰¹ our intuition.

4.1.3.2 Mean and variance of the intercept

⁴⁸⁰² Instead of estimating the intercept using data y_i following formula (4.20), we define
⁴⁸⁰⁴ the intercept as a random variable using the same formula but with random variable
⁴⁸⁰⁵ Y :

$$A = \bar{Y} - B\bar{x}. \tag{4.57}$$

4806 Below we show that $E[A] = a$:

$$\begin{aligned} E[A] &= E[\bar{Y}] - E[B]\bar{x} \\ &= \frac{\sum_{i=1}^n Y_i}{n} - b \frac{\sum_{i=1}^n x_i}{n} \\ &= \frac{\sum_{i=1}^n (Y_i - bx_i)}{n} \\ &= \frac{\sum_{i=1}^n a}{n} \\ &= a. \end{aligned} \tag{4.58}$$

4807 The variance of A is equal to

$$\text{Var}[A] = \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{\sigma_x^2} \right), \tag{4.59}$$

4808 where

$$\sigma_x^2 = \frac{S_{xx}}{n} \tag{4.60}$$

4809 may be considered the estimated variance of the x_i data. The result can be
4810 derived as follows:

$$\begin{aligned} \text{Var}[A] &= \text{Var}[\bar{Y}] + \text{Var}[B]\bar{x}^2 \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2}{S_{xx}} \bar{x}^2 \\ &= \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{S_{xx}/n} \right). \end{aligned} \tag{4.61}$$

4811 This expression implies that as n goes to infinity, the variance $\text{Var}[A]$ also goes
4812 to zero. Thus, the intercept estimate is better when having more data points.

4813 4.1.3.3 Confidence intervals of slope and intercept

4814 The $(1 - \alpha) \times 100\%$ confidence interval of the slope is

$$\left(\hat{b} - t_{\alpha/2, n-2} \frac{s}{\sqrt{S_{xx}}}, \quad \hat{b} + t_{\alpha/2, n-2} \frac{s}{\sqrt{S_{xx}}} \right), \tag{4.62}$$

4815 where

$$s = \sqrt{\frac{SSE}{n-2}} \tag{4.63}$$

4816 is the unbiased estimator of σ , and the degrees of freedom (dof) of the t-distribution
4817 is $n - 2$.

4818 If the confidence interval does not include zero, then the slope is significantly
4819 different from zero at the $\alpha \times 100\%$ significance level. The Colorado TLR example
4820 has $\hat{b} = -0.006982$. For the 95% confidence interval, the R's quantile function `qt` of
4821 the t-distribution can produce the quantile interval corresponding to probabilities
4822 0.025 and 0.975:

```
4823 qt(c(.025, .975), df=22)
4824 #[1] -2.073873 2.073873
```

4825 The estimated standard error for \hat{b} is

$$\text{SE}[\hat{b}] = \frac{s}{\sqrt{S_{xx}}} = 0.0002913 \quad (4.64)$$

4826 based on the R output from `summary(reg)`. Thus, the confidence interval for \hat{b} is

$$\begin{aligned} & (-0.0069818 - 2.073873 \times 0.0002913, -0.0069818 + 2.073873 \times 0.0002913) \\ & = (-0.0076, -0.0064). \end{aligned} \quad (4.65)$$

4827 This confidence interval does not include zero, and hence the slope is significantly
4828 less zero at the 5% confidence level. The 95% confidence interval for TLR is (6.4,
4829 7.6) °C/km based on the Colorado data used in this chapter.

4830 The $(1 - \alpha) \times 100\%$ confidence interval of the intercept is

$$\left(\hat{a} - t_{\alpha/2, n-2} s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}, \hat{a} + t_{\alpha/2, n-2} s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} \right). \quad (4.66)$$

4831 The R output from `summary(reg)` gives the standard error of \hat{a}

$$\text{SE}[\hat{a}] = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} = 0.5460279 \quad (4.67)$$

4832 Thus, the confidence interval for \hat{a} is

$$\begin{aligned} & (33.4762157 - 2.073873 \times 0.5460279, 33.4762157 + 2.073873 \times 0.5460279) \\ & = (32.3438, 34.6086). \end{aligned} \quad (4.68)$$

4833 The mathematical expressions of the confidence interval can be derived as follows.
4834 Formulas (4.51), (4.55) (4.58) and (4.59) imply that

$$A \sim N \left(a, \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{\sigma_x^2} \right) \right), \quad (4.69)$$

$$B \sim N \left(b, \frac{\sigma^2}{S_{xx}} \right) \quad (4.70)$$

$$\sigma_x^2 = S_{xx}/n \quad (4.71)$$

4835 for a given σ^2 .

4836 When σ^2 is unknown and is to be estimated by s^2 , then A and B follow a t-distribution.
4837 To be exact, the t-statistics with $n - 2$ dof can be defined as

$$T_a = \frac{\hat{a} - a}{s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}}, \quad (4.72)$$

$$T_b = \frac{\hat{b} - b}{s / \sqrt{S_{xx}}}, \quad (4.73)$$

4838 where a and b are given, and \hat{a} and \hat{b} are estimated from data. The confidence

4839 interval $(-t_{\alpha/2,n-2}, t_{\alpha/2,n-2})$ for T_a and T_b leads to the confidence intervals for \hat{a}
4840 and \hat{b} . Thus, these two statistics T_a and T_b can be used for a hypothesis test.

4.1.3.4 Hypothesis test for the slope using the t-test

4841 We wish to test the hypothesis whether the slope is equal to a given constant β .
4842 The null and alternative hypotheses are

$$H_0 : b = \beta, \quad (4.74)$$

$$H_1 : b \neq \beta. \quad (4.75)$$

4843 We define a t-statistic

$$T_b = \frac{\hat{b} - \beta}{s/\sqrt{S_{xx}}}. \quad (4.76)$$

4844 At the $\alpha \times 100\%$ significance level, we reject H_0 if

$$|T_b| > t_{1-\alpha/2,n-2}. \quad (4.77)$$

4845 For the Colorado July TLR example, if we want to test whether the TLR is equal
4846 to $7.3 \text{ } ^\circ\text{C/km}$, we calculate

$$T_b = \frac{-0.0069818 - (-0.0073)}{0.0002913} = 1.092345. \quad (4.78)$$

4847 The quantile

$$t_{0.975,22} = 2.073873 \quad (4.79)$$

4848 Thus, $|T_b| < t_{0.975,22}$, and hence the null hypothesis is not rejected at the 5%
4849 significance level. Namely, the obtained TLR $7.0 \text{ } ^\circ\text{C/km}$ is not significantly different
4850 from the given value $7.3 \text{ } ^\circ\text{C/km}$.

4.1.3.5 Confidence intervals of prediction: The fitted model and the response variable

4851 For an assigned point x^* , which is not on the data points x_i , the fitted model
4852 prediction is

$$\hat{Y} = \hat{a} + \hat{b}x^*. \quad (4.80)$$

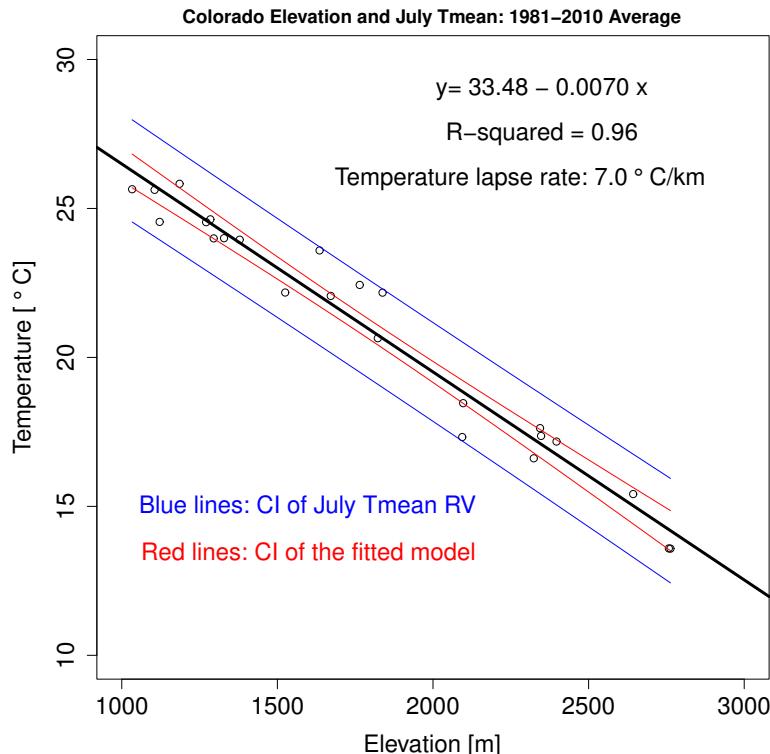
4853 The variance of this model prediction is

$$\text{Var}[\hat{a} + \hat{b}x^*] = \sigma^2 \left[\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right]. \quad (4.81)$$

4854 Thus, the $(1 - \alpha) \times 100\%$ confidence interval for $\hat{a} + \hat{b}x^*$ is given by

$$\hat{a} + \hat{b}x^* \pm t_{1-\alpha/2,n-2}s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}. \quad (4.82)$$

4858 This interval can be computed for any $x = x^*$ and is shown by the red lines in Fig.
 4859 4.3. It means that 95% of the chance for the model prediction to be between the
 4860 red lines. Note that the red lines diverge as x^* deviates from \bar{x} . This implies that
 4861 the confidence interval becomes larger as we go out of the range of the data.

**Fig. 4.3**

The confidence interval of the fitted model based on the Colorado July mean temperature data and their corresponding station elevations (Red lines). The confidence interval of the response variable Tmean (Blue lines).

4862 The error of the prediction at x^* is

$$e^* = Y - \hat{Y} = Y - (\hat{a} + \hat{b}x^*) \quad (4.83)$$

4863 is normally distributed if all the assumptions for the linear model are satisfied. The
 4864 error's mean is zero and variance is

$$\begin{aligned} \text{Var}[e^*] &= \text{Var}[Y - \hat{a} + \hat{b}x^*] \\ &= \text{Var}[Y] + \text{Var}[\hat{a} + \hat{b}x^*] \\ &= \sigma^2 + \sigma^2 \left[\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right]. \end{aligned} \quad (4.84)$$

4865 This result implies that the confidence interval for the response variable Y at an

4866 assigned x^* is

$$\hat{a} + \hat{b}x^* \pm t_{1-\alpha/2, n-2}s\sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}. \quad (4.85)$$

4867 The confidence interval is shown by the blue lines in Fig. 2.2. It means that 95% of
 4868 the chance for the Y values to be in between the blue lines. As expected, almost all
 4869 the Colorado TLR observed data lie between the blue lines. Only 5% of chance a
 4870 data point would lie outside the blue lines. According to Eq. (4.85), the blue lines
 4871 also diverge, but very slowly because

$$\frac{(x^* - \bar{x})^2}{S_{xx}} < 0.1310$$

4872 is much smaller than $1 + \frac{1}{n} = 1.0417$.

4873 Figure 4.3 can be produced by the following computer code.

```
#R plot Fig. 4.3: Confidence intervals of a regression model
4874 setwd("/Users/sshen/climstats")
4875 #Confidence interval of the linear model
4876 x1 = seq(max(x), min(x), len=100)
4877 n = 24
4878 xbar = mean(x)
4879 reg = lm(y ~ x)
4880 SSE = sum((reg$residuals)^2)
4881 s_squared = SSE/(length(y)-2)
4882 s = sqrt(s_squared)
4883 modTLR = 33.476216 + -0.006982*x1
4884 xbar = mean(x)
4885 Sxx = sum((x-xbar)^2)
4886 CIupperModel = modTLR +
4887   qt(.975, df=n-2)*s*sqrt((1/n)+(x1-xbar)^2/Sxx)
4888 CIlowerModel = modTLR -
4889   qt(.975, df=n-2)*s*sqrt((1/n)+(x1-xbar)^2/Sxx)
4890 CIupperResponse = modTLR +
4891   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x1-xbar)^2/Sxx)
4892 CIlowerResponse = modTLR -
4893   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x1-xbar)^2/Sxx)
4894
4895 setEPS() #Plot the figure and save the file
4896 postscript("fig0403.eps", height = 8, width = 8)
4897 par(mar=c(4.5,4.5,2.0,0.5))
4898 plot(x,y,
4899   ylim=c(10,30), xlim=c(1000,3000),
4900   xlab="Elevation [m]",
4901   ylab=bquote("Temperature["~degree~"C]"),
4902   main="Colorado[Elevation]and[July]Tmean:[1981-2010]Average",
4903   cex.lab=1.5, cex.axis=1.5)
4904 lines(x1,CIupperModel,type="l",col='red')
4905 lines(x1,CIlowerModel,type="l",col='red')
4906 lines(x1,CIupperResponse,type="l",col='blue')
4907 lines(x1,CIlowerResponse,type="l",col='blue')
4908 abline(reg,lwd=3)
4909 text(2280, 26,
4910 bquote("Temperature[lapse]rate:[7.0]~degree~"C/km)),
4911 cex=1.5)
```

```

4913 text(2350, 27.5,"R-squared=0.96", cex=1.5)
4914 text(2350, 29, "y=33.48+0.0070*x", cex=1.5)
4915 text(1600, 15, "Blue lines: CI of July Tmean RV",
4916   col="blue", cex=1.5)
4917 text(1600, 13.5,"Red lines: CI of the fitted model",
4918   col="red", cex=1.5)
4919 dev.off()

#Python plot Fig. 4.3: Confidence intervals of regression
reg1 = np.array(np.polyfit(x, y, 1))#regression
abline = reg1[1] + x*reg1[0]
# get confidence intervals
yl,yu,xd,rl,ru = linregress_CIs(x,y)
# plot the figure
fig, ax = plt.subplots(figsize=(13,12))
ax.plot(x,y, 'ko');
ax.plot(x, abline, 'k-');
ax.plot(xd, yu, 'r-')
ax.plot(xd, yl, 'r-')
ax.plot(xd, ru, 'b-')
ax.plot(xd, rl, 'b-')
ax.set_title("Colorado Elevation vs. July Tmean\n\n1981-2010 Average",size=25, pad = 20,
            fontweight = 'bold')
ax.set_xlabel("Elevation [m]", size = 25, labelpad = 20)
ax.set_ylabel(r"Temperature [$\degree C"]",
            size = 25, labelpad = 20);
ax.tick_params(length=6, width=2, labelsize=20);
ax.set_xticks(np.round(np.linspace(1000, 3000, 5), 2))
ax.set_yticks(np.round(np.linspace(10, 30, 5), 2))
ax.text(1600, 25, r"Temp lapse rate: 7.0$\degree C/km",
        color= 'k', size = 20)
ax.text(1750, 27,
        r"$y=%1.2f%1.4fx$%(reg1[1], reg1[0]),
        size = 20)
ax.text(1800, 26, r"$R-squared=%2f$" % Rsqu,
        color= 'k', size = 20)
ax.text(1050, 15, "Blue lines: CI of July Tmean RV",
        color= 'b', size = 20)
ax.text(1050, 14, "Red lines: CI of the fitted model",
        color= 'r', size = 20)
plt.show()

```

4920

4.1.3.6 When the assumptions for the linear regression model are violated

The computing and plotting of a linear regression is quite straightforward, and have been used widely in data analysis. However, users often are unaware of the assumptions behind the computing and results. You may wish to know what are the assumptions and how they are violated.

(a) Review of the assumptions of linear regression

4928 The linear regression model and its inference have four major assumptions:

- 4929 (i) Approximate linearity: The data support an approximate linear relationship
4930 between x and Y ;
- 4931 (ii) Normality: Both the model error term ϵ and Y are normally distributed;
- 4932 (iii) Constant variance: ϵ and Y have variance σ^2 , which does not change with
4933 respect to x ;
- 4934 (iv) Independence of the error term: $\text{Cor}(\epsilon_i, \epsilon_j) = \delta_{ij}$.

4935 Although the derivation of the estimation formulas for \hat{a} and \hat{b} does not need any
4936 of the above assumptions, the interpretation and inferences on \hat{a} , \hat{b} , $\hat{a} + \hat{b}x$, and
4937 other quantities are often based on these assumptions. If one or more assumptions
4938 are violated, the inference may not work. There are a variety of checking methods
4939 and resolution methods in literature, such as using a nonlinear relationship (e.g., a
4940 polynomial fitting model). R has a package `lmtest` to test linear regression models.
4941 A residual scatter plot (x_i, e_i) may be a first good step to intuitively identify possible
4942 violations of the assumptions (See Fig. 4.4). Ideally, the plot shows uniform noise for
4943 a good model. Figure 4.4 shows the residuals of the Colorado temperature regression
4944 against elevation, and seems showing this kind of noise without an obvious pattern.
4945 We may safely conclude that the assumptions of linearity and constant variance are
4946 satisfied.

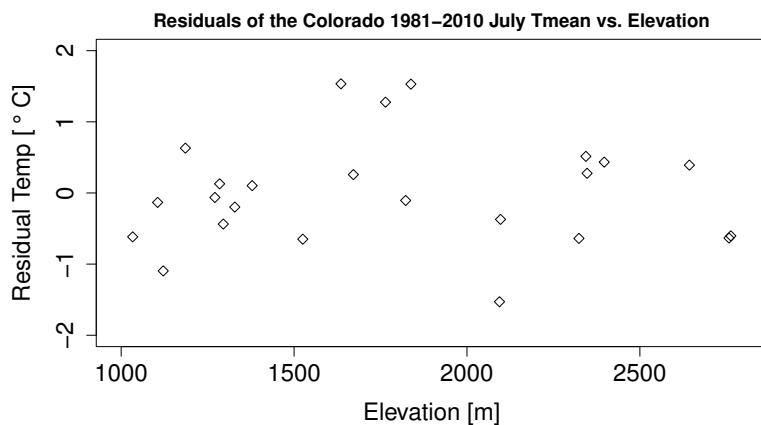


Fig. 4.4 Scatter plot of the residuals of the linear regression for the Colorado July mean temperature data against the elevations. The data are from the 24 USHCN stations in Colorado, and the July Tmean is the 1981-2010 average temperature.

4947 To be rigorous on the statistical inferences about the results of a linear regres-
4948 sion, a comprehensive residual analysis may be made to assess whether the studied
4949 dataset satisfies the linear regression assumptions. If some assumptions are violated
4950 for a given dataset, one may seek other methods to remediate the problems. One
4951 way is to make a data transformation or to use a new regression function so that
4952 the assumptions are satisfied for the transformed data or the new functional model.

Another way is to use a non-parametric method, which does not assume any distribution. The comprehensive residual analysis methods and the corresponding R and Python codes can be found online. This book is limited to the discussion of a few examples.

Violation of independence assumption (iv) often occurs in cases of y being a time series. In this case there could be correlations from one time to later times. This serial correlation effect leads to fewer independent time intervals, i. e., the degrees of freedom (dof) can be reduced, and the statistical inference needs special attention.

```
#R plot Fig. 4.4: Regression residuals
reg = lm(y~x)
setEPS() #Plot the figure and save the file
postscript("fig0404.eps", height = 4.5, width = 8)
par(mar=c(4.5,4.5,2.0,0.5))
plot(x, reg$residuals, pch=5,
      ylim=c(-2,2), xlim=c(1000,2800),
      xlab="Elevation [m]",
      ylab=bquote("Residual[Temp] [degree C]"),
      main="Residuals of the Colorado July Tmean vs. Elevation",
      cex.lab=1.5, cex.axis=1.5, cex.main = 1.2)
dev.off()
```

```
# Python plot Fig. 4.4: Regression residuals
reg1 = np.array(np.polyfit(x, y, 1))#regression
abline = reg1[1] + x*reg1[0]
# calculate residuals
r = y - abline
fig, ax = plt.subplots(figsize=(12,8))
ax.plot(x, r, 'kd')
ax.set_title("Residuals of the Colorado 1981-2010 July\nTmean vs. Elevation", fontweight = 'bold', size=25, pad = 20)
ax.set_xlabel("Elevation [$m$]", size = 25, labelpad = 20)
ax.set_ylabel(r"Residual[Temp] [$\degree C$]", size = 25, labelpad = 20);
ax.tick_params(length=6, width=2, labelsize=20);
ax.set_xticks(np.linspace(1000, 3000, 5))
ax.set_yticks(np.linspace(-2,2,5))
plt.show()
```

4973

(b) Test for normality

The normality and independence assumptions may not show in the residual scatter plot. You may use a Q-Q plot to test the normality, and Durbin-Watson test to check the independence. In the Colorado TLR example, the Q-Q plot is shown in Fig. 4.5. The residual data quantile points overlap fairly well with the theoretical line for the normal distribution. We may conclude that the data satisfy the normality assumption.

However, the QQ-plot is only a subjective visual test. A quantitative test may be used, such as the Kolmogorov-Smirnov (KS) test or Shapiro-Wilk (SW) test. The KS-test statistic and its p-value for the Colorado TLR data are $D = 0.20833$

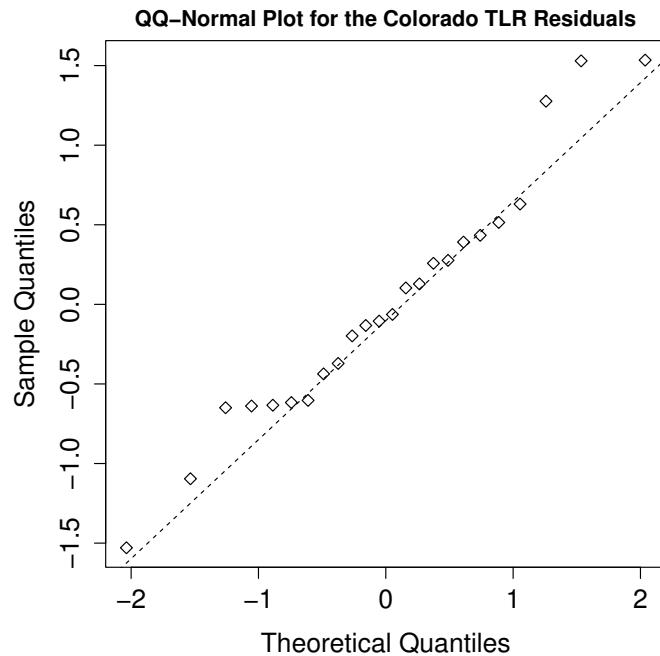


Fig. 4.5 Q-Q plot of the linear regression residuals for Colorado July mean temperature data at the USHCN 24 stations. The diagonal straight line is the theoretical line for the standard normal distribution.

and p-value = 0.686. The null hypothesis is not rejected, i.e., the temperature residuals from the linear regression are not significantly different from the normal distribution. The normality assumption is valid.

Figure 4.5 can generated by the following computer code.

```

4988 #R plot Fig. 4.5: Q-Q plot of quantiles
4989 reg = lm(y~x)
4990 setEPS() #Plot the figure and save the file
4991 postscript("fig0405.eps", height = 6, width = 6)
4992 par(mar=c(4.5,4.5,2.0,0.5))
4993 qqnorm(reg$residuals, pch=5,
4994   main="QQ-Normal Plot for the Colorado TLR Residuals",
4995   cex.lab = 1.4, cex.axis = 1.4)
4996 qqline(reg$residuals, lty=2)
4997 dev.off()

```

```

#Python plot Fig. 4.5: Q-Q plot of quantiles
fig, ax = plt.subplots(figsize=(12,8))#fig setup
reg1 = np.array(np.polyfit(x, y, 1))#regression
abline = reg1[1] + x*reg1[0]
r = y - abline# calculate residuals
#Q-Q plot as a probability plot: quantiles vs quantiles
pp1 = scistats.probplot(r, dist="norm", plot=ax,)
ax.set_title("QQ-Normal Plot for Colorado TLR Residuals",
            fontweight = 'bold', size = 25, pad = 20);
ax.set_ylabel("Sample Quantiles", size = 25, labelpad = 20);
ax.set_xlabel("Theoretical Quantiles",
              size = 25, labelpad = 20);
ax.tick_params(length=6, width=2, labelsize=20);
ax.set_xticks(np.round(np.linspace(-2, 2, 5), 2))
plt.show()

```

4998

4999 (c) Serial correlation

5000 The R package `lmtest` has a function to test for the independence of the error
 5001 term, using the Durbin-Watson (DW) test statistic defined as

$$DW = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}. \quad (4.86)$$

5002 The DW test is meant for a time series. We thus sort the `Tmean` data according to
 5003 the ascending order of elevation, and then compute DW.

5004 When $\text{cor}(e_i, e_j) = \sigma^2 \delta_{ij}$, we have $DW \approx 2$. For the Colorado July TLR data,
 5005 $DW = 2.3072$ and $p-value = 0.7062$, which can be computed by the following
 5006 code:

```

5007 #R code for the DW-test for independence
5008 #install.packages("lmtest")
5009 library(lmtest)
5010 ElevTemp=cbind(x,y, 1:24)
5011 #Sort the data for ascending elevation
5012 ElevTemp=ElevTemp[order(x),]
5013 reg1=lm(ElevTemp[,2] ~ ElevTemp[,1])
5014 dwtest(reg1)
5015 #DW = 2.3072, p-value = 0.7062

```

```

#Python code for the DW-test for independence
from statsmodels.stats.stattools import durbin_watson
orderET = np.arange(1,25)
ElevTemp = np.stack((x,y, orderET), axis = -1)
sorted_ind = np.argsort(ElevTemp[:,0], kind='mergesort')
dat1 = ElevTemp[sorted_ind]
# use the first order polynomial fit for linear regression
reg1 = np.array(np.polyfit(dat1[:,0], dat1[:,1], 1))
abline = reg1[1] + dat1[:,0]*reg1[0]
r = dat1[:,1] - abline # calculate residuals
#perform Durbin-Watson test
durbin_watson(r)
#2.307190038542777

```

5016

5017 The large p-value 0.7062 implies that the null hypothesis of no serial correlation
 5018 (i.e., independence) is not rejected, because DW 2.3072 is not too far away from 2,
 5019 which indicates independence of the error terms. Some literatures use $1.5 < DW <$
 5020 2.5 to conclude independence without checking the p-value.

5021 When the independence assumption is violated, the data have serial correlation,
 5022 which reduces dof and hence enlarges the confidence interval for the regression
 5023 results and makes the results not as reliable.

5024 (d) Non-parametric trend inference and heteroskedasty

5025 Even when some of the four assumptions for the linear regression model are
 5026 invalid, one can still estimate \hat{a} and \hat{b} and the linear model $\hat{a} + \hat{b}x$ in the same
 5027 way and can still interpret the study subject using the domain knowledge. In this
 5028 case, the statistical inference based on the normality assumption may not make
 5029 any sense, and the confidence interval calculations based on the t-distribution are
 5030 generally invalid. Now, you may use a non-parametric test to make an inference
 5031 on the residuals. For example, a non-parametric test for the trend can be applied,
 5032 such as Mann-Kendall (MK) trend test, and Theil-Sen's trend estimation and test.
 5033 The trend, if presenting, can be linear or nonlinear. These tests are meant for time
 5034 series. Thus, before applying these tests, the data should be sorted according to the
 5035 ascending order of the explainable variable. For the Colorado TLR example, the
 5036 MK-test for the linear regression residual can be computed by the following code.

```

#R code for the Mann-Kendall test
#install.packages("trend")
library(trend)
ElevTemp=cbind(x, y, 1:24)
#Sort the data for ascending elevation
ElevTemp=ElevTemp[order(x),]
reg1=lm(ElevTemp[,2] ~ ElevTemp[,1])
ElevTemp[,3]=reg1$residuals
mk.test(ElevTemp[,3])
#data: ElevTemp[, 3]
#z = 0.47128, n = 24, p-value = 0.6374
mk.test(ElevTemp[,2])
#z = -5.9779, n = 24, p-value = 2.261e-09

```

```

#Python code for the Mann-Kendall test
import pymannkendall as mk
orderET = np.arange(1,25)
ElevTemp = np.stack((x,y, orderET), axis = -1)
sorted_ind = np.argsort(ElevTemp[:,0],kind='mergesort')
dat1 = ElevTemp[sorted_ind]
# use the first order polynomial fit for linear regression
reg1 = np.array(np.polyfit(dat1[:,0], dat1[:,1], 1))
abline = reg1[1] + dat1[:,0]*reg1[0]
r = dat1[:,1] - abline# calculate residuals
#perform Durbin-Watson test
dat1[:,2] = r
print(mk.original_test(dat1[:,2]))#test for residual trend
#p=0.6374381847429542, z=0.47128365713220055
print(mk.original_test(dat1[:,1]))#test for temp trend
#p=2.260863496417187e-09, z=-5.97786112467686

```

5050

5051 The large p-value 0.6374 of the MK-test for residuals implies that the residuals
 5052 have no trend, linear or nonlinear. The small p-value 2.261e-09 of the MK-test for
 5053 the sorted temperature implies that the sorted July Tmean according to elevation
 5054 has a significant trend.

5055 The linear trend is frequently used to explain the global average temperature. The
 5056 temperature data with respect to time is clearly nonlinear. When the scatter plot
 5057 shows a clear pattern, then one or more of the four assumptions are usually violated.
 5058 When that happens, you may use different ways to remediate this situation. For
 5059 example, when linearity is violated, the problem may be intrinsically nonlinear and
 5060 a polynomial model may be a better fit than a linear model. When the constant
 5061 variance is violated, known as heteroskedastic data in statistics literature, then the
 5062 data standardization (i.e., the anomaly data being divided by the data standard
 5063 deviation) or logarithmic transform for the positive-valued data, may be used to
 5064 transform the data and to make the data homoskedastic, i.e., constant variance.

5065

4.2 Multiple linear regression

5066

5067 This section uses the method of matrices, which will be described in Chapter 5. If
 5068 you are not familiar with matrices, you may come back to this section after reading
 5069 the first two sections of Chapter 5.

5070

4.2.1 Calculating the Colorado TLR when taking location coordinates into account

5071

5072

5073 In the previous simple linear regression, the July temperature was assumed to
 5074 linearly depend only on the vertical coordinate of the station: elevation. However,

the temperature may also depend on the horizontal coordinates of the station: latitude and longitude. This subsection deals with this kind of problem of more than one explanatory variable. A multivariate linear regression model (or called multiple linear regression model) can be expressed as follows:

$$Y = b_0 + b_1x_1 + b_2x_2 + \cdots + b_mx_m + \epsilon, \quad (4.87)$$

where x_1, x_2, \dots, x_m are m explanatory variables, which are non-random and deterministic variables; ϵ is the model error with zero mean and variance σ^2 being a constant; Y is the response variable, which is random and has its variance equal to σ^2 , and expected value equal to

$$\mathbb{E}[Y] = b_0 + b_1x_1 + b_2x_2 + \cdots + b_mx_m; \quad (4.88)$$

and $b_0, b_1, b_2, \dots, b_m$ are parameters to be estimated from data $(x_{ij}, y_j), i = 1, 2, \dots, m$, and $j = 1, 2, \dots, n$.

In the Colorado July TLR example, we have x_1 as latitude, x_2 as longitude, x_3 as elevation, Y as the July air temperature, $m = 3$, and $n = 24$. The latitude and longitude data are as follows:

```
5088 lat=c(
5089 39.9919, 38.4600, 39.2203, 38.8236, 39.2425, 37.6742,
5090 39.6261, 38.4775, 40.6147, 40.2600, 39.1653, 38.5258,
5091 37.7717, 38.0494, 38.0936, 38.0636, 37.1742, 38.4858,
5092 8.0392, 38.0858, 40.4883, 37.9492, 37.1786, 40.0583
5093 )
5094 lon=c(
5095 -105.2667, -105.2256, -105.2783, -102.3486, -107.9631, -106.3247,
5096 -106.0353, -102.7808, -105.1314, -103.8156, -108.7331, -106.9675,
5097 -107.1097, -102.1236, -102.6306, -103.2153, -105.9392, -107.8792,
5098 -103.6933, -106.1444, -106.8233, -107.8733, -104.4869, -102.2189
5099 )
```

The elevation and temperature data are the same as those at the beginning of this chapter.

The computer code for the multiple linear regression of three variables is as follows.

```
5104 #R code for the TLR multivariate linear regression
5105 elev = x; temp = y #The x and y data were entered earlier
5106 dat = cbind(lat, lon, elev, temp)
5107 datdf = data.frame(dat)
5108 datdf[1:2,] #Show the data of the first two stations
5109 #      lat      lon    elev    temp
5110 # 39.9919 -105.2667 1671.5 22.064
5111 # 38.4600 -105.2256 1635.6 23.591
5112
5113 #Multivariate linear regression
5114 reg=lm(temp ~ lat + lon + elev, data = datdf)
```

```

5115 summary(reg) #Display the regression results
5116 #             Estimate Std. Error t value Pr(>|t|)
5117 #(Intercept) 36.4399561 9.4355746 3.862 0.000971 ***
5118 #   lat        -0.4925051 0.1320096 -3.731 0.001319 **
5119 #   lon         -0.1630799 0.0889159 -1.834 0.081564 .
5120 #   elev        -0.0075693 0.0003298 -22.953 7.67e-16 ***
5121 #Residual standard error: 0.6176 on 20 degrees of freedom
5122 #Multiple R-squared:  0.979

```

```

#Python code for the TLR multivariate linear regression
from sklearn import linear_model
lat=np.array([
39.9919, 38.4600, 39.2203, 38.8236, 39.2425, 37.6742,
39.6261, 38.4775, 40.6147, 40.2600, 39.1653, 38.5258,
37.7717, 38.0494, 38.0936, 38.0636, 37.1742, 38.4858,
8.0392, 38.0858, 40.4883, 37.9492, 37.1786, 40.0583
])
lon=np.array([
-105.2667,-105.2256,-105.2783,-102.3486,-107.9631,-106.3247,
-106.0353,-102.7808,-105.1314,-103.8156,-108.7331,-106.9675,
-107.1097,-102.1236,-102.6306,-103.2153,-105.9392,-107.8792,
-103.6933,-106.1444,-106.8233,-107.8733,-104.4869,-102.2189
])
elev = x; temp = y #The x and y data were entered earlier
dat = np.stack((lat, lon, elev), axis = -1)
print(dat[0:2,:]) #Show the data of the first two stations
#[[ 39.9919 -105.2667 1671.5 ]
#[ 38.46 -105.2256 1635.6 ]]
#Multivariate linear regression
xdat = dat
ydat = temp
regr = linear_model.LinearRegression()
multi_reg = regr.fit(xdat, ydat)
print('Intercept:\n', regr.intercept_)
print('Coefficients:\n', regr.coef_)

```

5123

According to this three-variable linear regression, the TLR is the regression coefficient for elevation, i.e., -0.0075694 . This means $7.6^{\circ}\text{C}/\text{km}$, larger than the one estimated earlier using only one variable: elevation.

The quantile of $t_{0.975,20} = 2.085963$. The dof is 20 now because four parameters have been estimated. The confidence interval of TLR can be computed below

$$\begin{aligned}
& (-0.0075694 - 2.085963 \times 0.0003298, -0.0075694 + 2.085963 \times 0.0003298) \\
& = (-0.008257351, -0.006881449).
\end{aligned} \tag{4.89}$$

5129 The confidence interval for TLR at 95% confidence level is $(6.9, 8.3)^{\circ}\text{C}/\text{km}$.

5130 The R^2 values is very large: 0.979, which means that the linear model data can 5131 explain 98% of the variance of observed temperature data.

5132 4.2.2 Formulas for estimating parameters in the multiple linear 5133 regression 5134

5135 The n groups of data for the explainable x_1, x_2, \dots, x_m and response variable Y
5136 are as follows:

$$x_{ij}, i = 1, 2, \dots, m, \text{ and } j = 1, 2, \dots, n; \quad (4.90)$$

5137 and

$$y_1, y_2, \dots, y_n \quad (4.91)$$

5138 The data and their corresponding regression coefficients are written in the matrix
5139 form as follows:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (4.92)$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{m1} \\ 1 & x_{12} & x_{22} & \cdots & x_{m2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{mn} \end{bmatrix} \quad (4.93)$$

$$\hat{\mathbf{b}} = \begin{bmatrix} \hat{b}_0 \\ \hat{b}_1 \\ \hat{b}_2 \\ \vdots \\ \hat{b}_n \end{bmatrix} \quad (4.94)$$

5140 The data, linear model prediction, and the corresponding residuals (i.e., the pre-
5141 diction errors of the linear model) can be written as follows:

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times (m+1)} \hat{\mathbf{b}}_{(m+1) \times 1} + \mathbf{e}_{n \times 1} \quad (4.95)$$

5142 where the residual vector is

$$\mathbf{e}_{n \times 1} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad (4.96)$$

5143 Multiplying Eq. (4.95) by \mathbf{X}^t from the left and enforcing

$$(\mathbf{X}^t)_{(m+1) \times n} \mathbf{e}_{n \times 1} = \mathbf{0}, \quad (4.97)$$

5144 you can obtain the estimate of the regression coefficients

$$\hat{\mathbf{b}}_{(m+1) \times 1} = [(\mathbf{X}^t)_{(m+1) \times n} \mathbf{X}_{n \times (m+1)}]^{-1} (\mathbf{X}^t)_{(m+1) \times n} \mathbf{y}_{n \times 1}. \quad (4.98)$$

5145 Condition (4.97) means the each column vector of the $\mathbf{X}_{n \times (m+1)}$ matrix is perpendicular
 5146 to the residual vector $\mathbf{e}_{n \times 1}$. For the first column, the condition corresponds
 5147 to the assumption of zero mean of the model error:

$$\sum_{i=1}^n e_i = 0. \quad (4.99)$$

5148 For the remaining columns, the condition means that the residual vector is perpen-
 5149 dicular to the data vector for each explanatory variable. This implies that the
 5150 residual vectors Euclidean distance to the data vectors are minimized, i.e., the min-
 5151 imizing sum of squared errors (SSE). Therefore, Equation (2.83) is the least square
 5152 estimate of the regression coefficients.

5153 As illustrated in the previous subsection, to implement the R estimate of the
 5154 regression coefficients, we put data \mathbf{X} and \mathbf{y} in a single matrix in the form of
 5155 `datdf=data.frame(cbind(X,y))` with proper column names, such as

5156 `colnames(datdf) <- c('x1', 'x2', 'x3', 'y')`

5157 for the case of three explanatory variables. Then, use the R command to make the
 5158 linear model estimate

5159 `reg=lm(y ~ x1 + x2 + x3, data = datdf)`

5160 Finally, `summary(reg)` outputs all the important regression results.

5161 The R command `reg$` allows to display the specific result, such as regression
 5162 coefficients

```
5163 round(reg$coefficients, digits=5)
5164 # (Intercept)      lat          lon          elev
5165 #     36.43996    -0.49251    -0.16308    -0.00757
```

5166 Thus, the multiple linear model for the Colorado July temperature is

$$\text{Temp}[\text{°C}] = 36.43527 - 0.49255 \times \text{Lat} - 0.16314 \times \text{Lon} - 0.00757 \times \text{Elev} [\text{m}]. \quad (4.100)$$

5167 The confidence interval for the linear model

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + \cdots + b_m x_m \quad (4.101)$$

5168 at a given point $\mathbf{x}^* = (1, x_1^*, x_2^*, \dots, x_m^*)$ is given by the following formula

$$\hat{\mathbf{b}}^t \mathbf{x}^* \pm t_{1-\alpha/2, n-m} s \sqrt{(\mathbf{x}^*)^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}^*}. \quad (4.102)$$

5169 The confidence interval for the response variable Y at a given point

$$\mathbf{x}^* = (1, x_1^*, x_2^*, \dots, x_m^*)$$

5170 is

$$\hat{\mathbf{b}}^t \mathbf{x}^* \pm t_{1-\alpha/2, n-m} s \sqrt{1 + (\mathbf{x}^*)^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}^*}, \quad (4.103)$$

5171 where s is the same as that defined earlier for the simple linear regression:

$$s = \sqrt{\frac{SSE}{n-2}} \quad (4.104)$$

5172 If $m = 1$, the above two formulas are the same as those for the confidence intervals
 5173 given by Eqs. (4.82) and (4.85) for the simple linear regression.

5174 **4.3 Nonlinear fittings using the multiple linear**
 5175 **regression**

5177 **4.3.1 Diagnostics of linear regression: An example of global**
 5178 **temperature**

5180 When the data have strong nonlinearity, the scatter plot of residuals will show
 5181 obvious patterns, such as that shown in Fig. 4.6 for the linear regression of the
 5182 global average annual mean surface air temperature anomalies from 1880 to 2018
 5183 with respect to the 1971-2000 climatology. The linearity assumption of the simple
 5184 linear regression is clearly violated. The independence assumption is also violated.

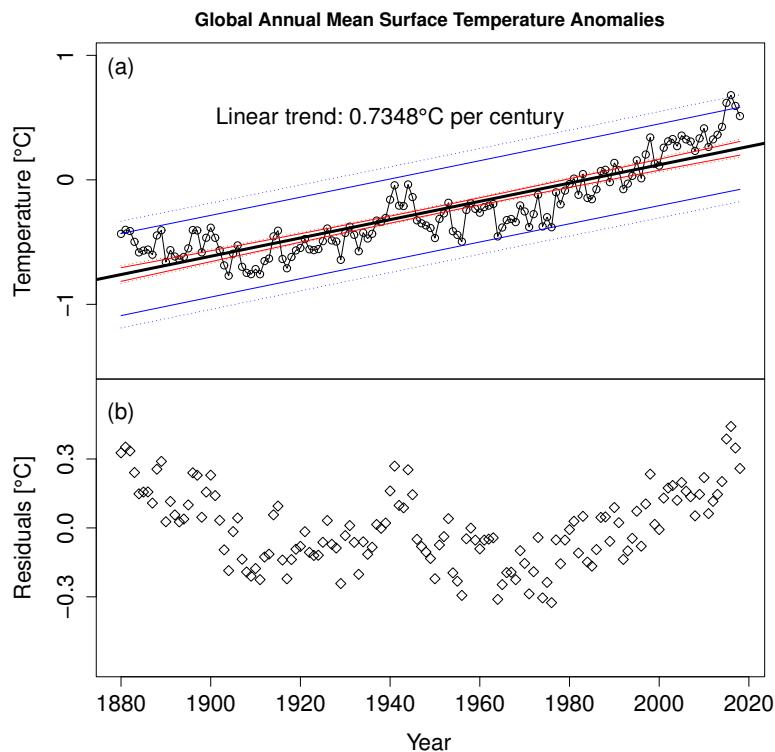


Fig. 4.6 (a) Linear regression of the global average annual mean land and ocean surface air temperature anomalies with respect to the 1971-2000 climatology based on the NOAAGlobalTemp dataset (Zhang et al. 2019); (b) Scatter plot of the linear regression residuals.

5185 The following computer code can plot Fig. 4.6 and make the diagnostics of the
 5186 linear regression.

5187 `#R plot Fig. 4.6: Regression diagnostics`

```

5188 setwd("/Users/sshen/climstats")
5189 dtmean<-read.table(
5190   "data/aravg.ann.land_ocean.90S.90N.v5.0.0.201909.txt",
5191   header=F)
5192 dim(dtmean)
5193 #[1] 140 6
5194 x = dtmean[1:139,1]
5195 y = dtmean[1:139,2]
5196 reg = lm(y ~ x) #linear regression
5197 reg
5198 #(Intercept) yrtime
5199 #-14.574841 0.007348
5200
5201 #Confidence interval of the linear model
5202 xbar = mean(x)
5203 SSE = sum((reg$residuals)^2)
5204 s_squared = SSE/(length(y)-2)
5205 s = sqrt(s_squared)
5206 modT = -14.574841 + 0.007348 *x
5207 xbar = mean(x)
5208 Sxx = sum((x-xbar)^2)
5209 n = length(y)
5210 CIupperModel= modT +
5211   qt(.975, df=n-2)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5212 CIlowerModel= modT -
5213   qt(.975, df=n-2)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5214 CIupperResponse= modT +
5215   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5216 CIlowerResponse= modT -
5217   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5218
5219 CIupperModelr= modT +
5220   qt(.975, df=5)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5221 CIlowerModelr= modT -
5222   qt(.975, df=5)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5223 CIupperResponser= modT +
5224   qt(.975, df=5)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5225 CIlowerResponser= modT -
5226   qt(.975, df=5)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5227
5228 setEPS() #Plot the figure and save the file
5229 postscript("fig0406.eps", height = 8, width = 8)
5230 par(mfrow=c(2,1))
5231 par(mar=c(0,4.5,2.5,0.7))
5232 plot(x, y, ylim = c(-1.5, 1),
5233   type="o", xaxt="n", yaxt="n",
5234   cex.lab=1.4, cex.axis=1.4,
5235   xlab="Year", ylab=bquote("Temperature["~degree~"C]"),
5236   main="Global[Annual]Mean[Surface]Temperature[Anomalies]",
5237   cex.lab=1.4, cex.axis=1.4
5238 )
5239 axis(side = 2, at = c(-1.0, 0, 1.0), cex.axis = 1.4)
5240 abline(reg, col="black", lwd=3)
5241 lines(x,CIupperModel,type="l",col='red')
5242 lines(x,CIlowerModel,type="l",col='red')
5243 lines(x,CIupperResponse,type="l",col='blue')

```

```
5244 lines(x,CIlowerResponse,type="l",col='blue')
5245
5246 lines(x,CIupperModelr,type="l", lty = 3, col='red')
5247 lines(x,CIlowerModelr,type="l", lty = 3, col='red')
5248 lines(x,CIupperResponser,type="l",lty = 3, col='blue')
5249 lines(x,CIlowerResponser,type="l",lty = 3, col='blue')
5250
5251 text(1940, 0.5,
5252     bquote("Linear_trend:~0.7348`^degree~"~Cuper~century`),
5253     col="black",cex=1.4)
5254 text(1880, 0.9, "(a)", cex=1.4)
5255 par(mar=c(4.5,4.5,0,0.7))
5256 plot(x, reg$residuals, ylim = c(-0.6,0.6),
5257     pch=5, cex.lab=1.4, cex.axis=1.4,
5258     yaxt = 'n', xlab="Year",
5259     ylab=bquote("Residuals["~degree~"~C]"))
5260 axis(side = 2, at = c(-0.3, 0, 0.3), cex.axis = 1.4)
5261 text(1880, 0.5, "(b)", cex=1.4)
5262 dev.off()
```

```

#Python plot Fig. 4.6: Regression diagnostics
# Change current working directory
os.chdir("/Users/sshenn/climstats/")
# Read the data
dtmean = np.array(read_table(
    "data/aravg.ann.land_ocean.90S.90N.v5.0.0.201909.txt",
    header = None, delimiter = "\s+"))
xDT = dtmean[0:139, 0]
yDT = dtmean[0:139, 1]
# Make a linear fit, i.e., linear regression
regLin = np.array(np.polyfit(xDT, yDT, 1))
ablinereg = regLin[1] + xDT*regLin[0]
yld,yud,xdd,rld,rud = linregress_CIs(xDT,yDT)

# Plot the figure
fig, ax = plt.subplots(2, 1, figsize=(12,12))
ax[0].plot(xDT, yDT, 'ko-')
ax[0].plot(xDT, ablinereg, 'k-')
ax[0].plot(xdd, yld, 'r-')
ax[0].plot(xdd, yud, 'r-')
ax[0].plot(xdd, rld, 'b-')
ax[0].plot(xdd, rud, 'b-')
ax[1].plot(xDT, yDT - ablinereg, 'kd')
ax[0].set_title(
    "Global\u2022Annual\u2022Mean\u2022Land\u2022and\u2022Ocean\u2022Surface\u2022\nTemperature\u2022Anomalies", fontweight = 'bold',
    size = 25, pad = 20)
ax[0].set_ylabel("Temperature\u2022[$\u00b0degree$C]", size = 25, labelpad = 20)
ax[0].tick_params(length=6, width=2, labelsize=20);
ax[0].set_yticks(np.round(np.linspace(-1, 1, 5), 2))
ax[0].text(1880, 0.4,
           "Linear\u2022temp\u2022trend\u20220.7348\u2022deg\u2022C\u2022per\u2022century",
           color = 'k', size = 20)
ax[0].text(1877, 0.6, "(a)", size = 20)
ax[0].axes.get_xaxis().set_visible(False)
ax[1].tick_params(length=6, width=2, labelsize=20);
ax[1].set_yticks(np.round(np.linspace(-0.5, 0.5, 5), 2))
ax[1].set_ylabel("Residuals\u2022[$\u00b0degree$C]", size = 25, labelpad = 20)
ax[1].text(1877, 0.4, "(b)", size = 20)
ax[1].set_xlabel("Year", size = 25, labelpad = 20)
fig.tight_layout(pad=-1.5)
plt.show()

```

5263

5264 Based on the visual check of Fig. 4.6(b), the constant variance assumption seems
 5265 satisfied. The KS-test shows that the normality assumption is also satisfied.

```

5266 #Kolmogorov-Smirnov (KS) test for normality
5267 library(fitdistrplus)
5268 resi_mean = mean(reg$residuals)
5269 resi_sd = sd(residuals)
5270 test_norm = rnorm(length(residuals),
                     mean = 0, sd = 1)
5271

```

```

5272 testvar = (reg$residuals - resi_mean)/resi_sd
5273 ks.test(testvar, test_norm)
5274 #D = 0.057554, p-value = 0.9754
5275 #The normality assumption is accepted
5276
5277 #Diagnostics on independence and normality
5278 # Durbin-Watson (DW) test for independence
5279 dwtest(reg)
5280 #DW = 0.45235, p-value < 2.2e-16
5281 #The independence assumption is rejected
5282
5283 #degrees of freedom and critical t values
5284 rho1 = acf(y)[[1]][2] #Auto-correlation function
5285 rho1 #[1] 0.9270817
5286 edof = (length(y) - 2)*(1 - rho1)/(1 + rho1)
5287 edof #[1] 5.183904 effective degrees of freedom
5288 qt(.975, df=137) #[1] 1.977431 critical t value
5289 qt(.975, df=5) #[1] 2.570582 critical t value

```

```

#Kolmogorov-Smirnov (KS) test for normality
import statistics
from scipy.stats import kstest
resi = yDT - ablinereg
testvar = (resi - np.mean(resi))/statistics.stdev(resi)
kstest(testvar, 'norm')#
#KstestResult(statistic=0.0751, pvalue=0.3971)
#The normality assumption is accepted

#perform Durbin-Watson test for independence
from statsmodels.stats.stattools import durbin_watson
durbin_watson(resi)
#0.45234915992769864 not in (1.5, 2.5)
#The independence assumption is rejected

#calculate autocorrelations of temp data for edof
import statsmodels.api as sm
autocorr = sm.tsa.acf(yDT)
rho1 = autocorr[1]
edof = (yDT.size - 2)*(1 - rho1)/(1 + rho1)
edof #[1] 5.183904 effective degrees of freedom

#Compare the critical t values with different edof
import scipy.stats
scipy.stats.t.ppf(q=.975, df=137)
##critical t value 1.9774312122928936
scipy.stats.t.ppf(q=.975, df=5)
##critical t value 2.5705818366147395

```

5290

5291 The DW-test shows that the independence assumption is violated. This implies
 5292 the existence of serial correlation, which in turn implies a smaller dof, and hence
 5293 larger $t_{1-\alpha, \text{dof}}$. Thus, the confidence intervals for model \hat{Y} are wider than the red
 5294 lines in Fig. 4.6(a), and those for the Y values are wider than the blue lines in

5295 the same figure. Therefore, the regression results are less reliable. In this case,
 5296 you need to compute the effective dof (edof), which is less than $n - 2$. For the
 5297 NOAAGlobalTemp time series, the one-year time lag serial correlation is $\rho_1 =$
 5298 0.9271, the edof under the assumption of an autoregression one process AR(1) is
 5299 approximately equal to

$$\text{edof} = \frac{1 - \rho_1}{1 + \rho_1} \times \text{dof}. \quad (4.105)$$

5300 Consequently, the effective dof is reduced to 5, compared to the non-serial corre-
 5301 lation case dof 137. R can compute `qt(.975, df=137)` and yield 1.977431, and
 5302 compute `qt(.975, df=5)` and yield 2.570582. Thus, the actual confidence interval
 5303 at the 95% confidence level for the linear regression is about 25% wider. Namely,
 5304 the dotted color lines are wider than the solid color lines.

5305 The W pattern of the residuals shown in Fig. 4.6(b) implies that the linearity
 5306 assumption is violated. To remediate the nonlinearity, we fit the data to a nonlinear
 5307 model and hope the residuals do not show clear patterns. We will use the third-order
 5308 polynomial to illustrate the procedure.

5309 5310 4.3.2 Fit a third order polynomial

5311 Figure 4.7(a) shows the third order polynomial fitting to the NOAAGlobalTemp
 5312 dataset:

$$y = b_0 + b_1x + b_2x^2 + b_3x^3 + \epsilon. \quad (4.106)$$

5313 This fitting can be estimated by the method of multiple linear regression with
 5314 three variables

$$x_1 = x, x_2 = x^2, x_3 = x^3. \quad (4.107)$$

5315 Then, a computer code for the multiple linear regression can be used to estimate
 5316 the coefficients of the polynomial and plot Fig. 4.7.

```
5317 #R plot Fig. 4.7: Polynomial fitting
5318 x1=x
5319 x2=x1^2
5320 x3=x1^3
5321 dat3=data.frame(cbind(x1,x2,x3,y))
5322 reg3 = lm(y ~ x1 + x2 + x3, data=dat3)
5323 # simply use
5324 # reg3 = lm(y ~ x + I(x^2) + I(x^3))
5325 setEPS() #Plot the figure and save the file
5326 postscript("fig0407.eps", height = 8, width = 8)
5327 par(mfrow=c(2,1))
5328 par(mar=c(0,4.5,2.5,0.7))
5329 plot(x, y, type="o", xaxt="n",
5330       cex.lab=1.4, cex.axis=1.4, xlab="Year",
5331       ylab=bquote("Temperature["~degree~"C]"),
5332       main="Global~Annual~Mean~Surface~Temperature~Anomalies",
5333       cex.lab=1.4, cex.axis=1.4)
5334 lines(x, predict(reg3), col="black", lwd=3)
5335 reg3
```

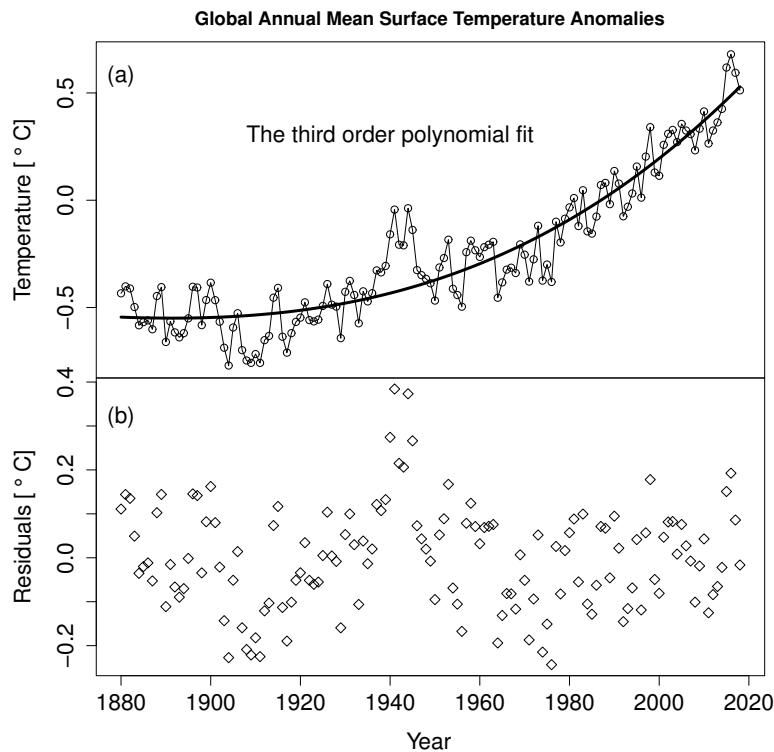


Fig. 4.7 (a) Fit a third order polynomial to the global average annual mean land and ocean surface air temperature anomalies with respect to the 1971-2000 climatology based on the NOAAGlobalTemp dataset (Zhang et al. 2019); (b) Scatter plot of the corresponding residuals.

```

5336  #(Intercept)           x1           x2           x3
5337  #-1.426e+03   2.333e+00   -1.271e-03   2.308e-07
5338  text(1940, 0.3,
5339    "The third order polynomial fit",
5340    col="black", cex=1.4)
5341  text(1880, 0.58, "(a)", cex=1.4)
5342  par(mar=c(4.5,4.5,0,0.7))
5343  plot(x1, reg3$residuals,
5344    pch=5, cex.lab=1.4, cex.axis=1.4,
5345    xlab="Year", ylab=bquote("Residuals["~degree~"C]"))
5346  text(1880, 0.32, "(b)", cex=1.4)
5347  dev.off()

```

```

os.chdir("/Users/sshenn/climstats/")
# Read the data
dtmean = np.array(read_table(
    "data/aravg.ann.land_ocean.90S.90N.v5.0.0.201909.txt",
    header = None, delimiter = "\s+"))
xDT = dtmean[0:139, 0]
yDT = dtmean[0:139, 1]
# Create trend line
reg3 = np.array(np.polyfit(xDT, yDT, 3))
ablineDT3 = reg3[3] + xDT*reg3[2] + \
            (xDT**2)*reg3[1] + (xDT**3)*reg3[0]
fig, ax = plt.subplots(2, 1, figsize=(12,12))
ax[0].plot(xDT, yDT, 'ko-')
ax[0].plot(xDT, ablineDT3, 'k-')

ax[1].plot(xDT, yDT - ablineDT3, 'kd')
ax[0].set_title("Global Annual Mean Land and Ocean\nSurface Temperature Anomalies",
                 fontweight = 'bold', size = 25, pad = 20)
ax[0].tick_params(length=6, width=2, labelsize=20);
ax[0].set_yticks(np.round(np.linspace(-1, 1, 5), 2))
ax[0].set_ylabel("Temperature [$\degree$C]",
                  size = 25, labelpad = 20)
ax[0].text(1900, 0.3, "The third-order polynomial fit",
           color = 'k', size = 20)
ax[0].text(1877, 0.6, "(a)", size = 20)
ax[0].axes.get_xaxis().set_visible(False)
ax[1].tick_params(length=6, width=2, labelsize=20);
ax[1].set_yticks(np.round(np.linspace(-0.4, 0.4, 5), 2))
ax[1].set_ylabel("Residuals [$\degree$C]",
                  size = 25, labelpad = 20)
ax[1].set_xlabel("Year", size = 25, labelpad = 20)
ax[1].text(1877, 0.35, "(b)", size = 20)
fig.tight_layout(pad=-1.5)
plt.show()

```

5348

5349 R has another simpler command for the third order polynomial regression:

5350 `lm(y ~ poly(x, 3, raw=TRUE))`

5351 where, `raw=TRUE` means using the raw polynomial form written like formula (4.106).
 5352 Another option is `raw=False`, which means that the data are fitted to an orthogonal
 5353 polynomial.

5354 Comparison of Figs. 4.6(b) and 4.7(b) shows that the W-shape pattern of the
 5355 residuals in the third-order polynomial fitting is weaker than that for the linear
 5356 regression. Nonetheless, the W-shape pattern is still clear. Figure 4.7(b) visually
 5357 suggests that the constant variance assumption is satisfied. The KS-test shows
 5358 that the normality assumption is also satisfied. However, the DW-test shows the
 5359 existence of serial correlation. The computing of these tests are left as exercise
 5360 problems.

5361 You can try to fit a ninth order orthogonal polynomial. This can eliminate the W-

5362 shape nonlinear pattern of the residuals. You can also try to fit other polynomials
 5363 or functions.

5364 4.4 Linear Regression by Weighted Least Squares

5366 The usual regression minimizes the mean square error:

$$\epsilon^2 = \sum_i (y_i - \hat{y}_i)^2. \quad (4.108)$$

5367 However, sometimes the square error may have different importance for different
 5368 observations i . The importance is quantified by a weight w_i . Thus, the weight least
 5369 square error is defined as

$$\varepsilon^2 = \sum_i w_i (y_i - \hat{y}_i)^2. \quad (4.109)$$

5370 The weights w_i are often required to be non-negative and sometimes to be normal-
 5371 ized

$$\sum_i w_i = 1, \quad w_i \geq 0, \quad i = 1, 2, 3, \dots \quad (4.110)$$

5372 As an example of weights, the gridded global climate data are often defined on
 5373 a uniform latitude-longitude grid. A grid box covers more area over the equatorial
 5374 region than a polar region. From the spherical coordinates, you can find that the
 5375 area-weight is $\cos \phi_i$, where ϕ_i is the latitude of the grid box's centroid.

5376 The weighted least squares regression can be converted into the usual least square
 5377 regression for the weighted field and data:

$$U_i = \sqrt{w_i} y_i, \quad \hat{U}_i = \sqrt{w_i} \hat{y}_i, \quad i = 1, 2, 3, \dots \quad (4.111)$$

5378 Here, $\sqrt{w_i}$ are often referred to as weight-factor. In climate science, $\sqrt{w_i} = \sqrt{\cos \phi_i}$
 5379 are called area-factor. In terms of the weighted data, the weighted least squares
 5380 becomes the usual least squares:

$$\varepsilon^2 = \sum_i (U_i - \hat{U}_i)^2. \quad (4.112)$$

5381 When considering 3D data in space, we may need to introduce the volume-factor
 5382 $\sqrt{V_i}$. When heat capacity, mass-density and other physical or chemical properties
 5383 in addition to geometry need to be considered, then a more complex weight should
 5384 be used. For more examples, see Bui et al. (2023), Lafarga et al. (2023)), and Shen
 5385 and North (2023).

5386

4.5 Chapter summary

5387

5388 This chapter has introduced three regression methods commonly used in climate
5389 science: (a) simple linear regression of a single variate, (b) multiple linear regression,
5390 and (c) nonlinear fitting. For (a) and (b) we used the data of surface air temperature
5391 and elevation of the 24 USHCN stations in the state of Colorado, USA. For (c) we
5392 used the NOAAGlobalTemp's global average annual mean temperature data from
5393 1880 to 2018.

5394 A linear regression model has four fundamental assumptions:

- 5395 (i) Approximate linearity between x and Y ,
- 5396 (ii) Normal distribution of ϵ and Y ,
- 5397 (iii) Constant variance of ϵ and Y , and
- 5398 (iv) Independence of error terms: $\text{Cor}(\epsilon_i, \epsilon_j) = \sigma^2 \delta_{ij}$.

5399 Assumptions (i) and (iii) can be verified by visually examining the scatter plot
5400 of residuals. Assumption (ii) may be verified by the KS-test. Assumption (iv) can
5401 be verified by the DW-test. When one or more assumptions are violated, you may
5402 improve your model, such as using a nonlinear model, or transforming the data.

References and Further Readings

- 5404 [1] Cordova, M., R. Celleri, C.J. Shellito, J. Orellana-Alvear, A. Abril, and G.
 5405 Carrillo-Rojas, 2016: Near-surface air temperature lapse rate over complex ter-
 5406 rain in the Southern Ecuadorian Andes: implications for temperature mapping.
 5407 *Arctic, Antarctic, and Alpine Research*, 48, 673 - 684.

Figure 3 of this paper shows the linear regression with $R^2 = 0.98$ for the temperature data of nine stations whose elevations are in the range from 2,610 to 4,200 meters. The regression implies a lapse rate of $6.88 \text{ }^\circ\text{C/km}$ for the annual mean temperature.

- 5409 [2] Fall, P.L., 1997: Timberline fluctuations and late Quaternary paleoclimates in
 5410 the Southern Rocky Mountains, Colorado. *Geological Society of America Bul-*
 5411 *letin*, 109, 1306-1320.

Figure 2a of this paper shows a lapse rate of $6.9 \text{ }^\circ\text{C/km}$ for the mean July temperature based on the data of 104 stations and a linear regression with $R^2 = 0.86$. Figure 2b shows the annual mean temperature TLR equal to $6.0 \text{ }^\circ\text{C/km}$ with $R^2 = 0.80$.

- 5413 [3] Graybill, F.A., H.K. Iyer, 1994: *Regression Analysis: Concepts and Applications*.
 5414 Duxbury Press, Belmont, California, 701pp.

This book clearly outlines the assumptions of regression. It focuses on concepts and applications about solving practical statistical problems.

- 5416 [4] Menne, M.J., C.N. Williams, and R.S. Vose, 2009: The United States Histori-
 5417 cal Climatology Network monthly temperature data Version 2. *Bulletin of the*
 5418 *American Meteorological Society*, 90, 993-1007.

This paper is one of the series of papers on the USHCN datasets prepared at the NOAA National Centers for Environmental Information and have been widely used since the early 1990s. The network includes 1,218 stations and has both monthly and daily data of temperature and precipitation publicly available. URL: <https://www.ncdc.noaa.gov/ushcn>

- 5420 [5] Zhang, H.-M., B. Huang, J. Lawrimore, M. Menne, T. M. Smith, 2019:
 5421 NOAA Global Surface Temperature Dataset (NOAAGlobalTemp), Version
 5422 5.0 (Time Series). NOAA National Centers for Environmental Information.
 5423 doi:10.7289/V5FN144H [Access date: March 2021].

5424 NOAA Merged Land Ocean Global Surface Temperature Analysis
 (NOAAGlobalTemp) dataset was produced and maintained by the
 NOAA National Centers for Environmental Information. This is monthly
 anomaly data with respect to the 1971-2000 climatology. It covers time
 from January 1880 to present, and has a spatial resolution $5^\circ \times 5^\circ$. The
 data can be visualized at www.4dvd.org.

5425

Exercises

- 5426
- 5427 **4.1** (a) Compute the temperature lapse rate for August using the TOB mean
 5428 temperature data from the 24 USHCN stations in Colorado during the period
 5429 of 1981-2010.
 5430 (b) Plot the figure similar to Fig. 4.1 for the data and results of Part (a).
- 5431 **4.2** Repeat the previous problem but for January. Compare your January results
 5432 with those of July in Fig. 4.1 and its related text.
- 5433 **4.3** (a) Compute the temperature lapse rate for the annual mean temperature
 5434 based on the 24 USHCN stations in Colorado during the period of 1981-2010.
 5435 (b) Plot the figure similar to Fig. 4.1 for the data and results of Part (a).
- 5436 **4.4** (a) Compute the annual mean temperature lapse rate for a high mountainous
 5437 region and for a period of your interest.
 5438 (b) Plot the figure similar to Fig. 4.1 for the data and results of Part (a).
- 5439 **4.5** Show that the orthogonality condition Eq. (4.27)

$$\mathbf{e} \cdot \mathbf{x}_a = 0 \quad (4.113)$$

5440 is equivalent to the condition of minimizing SSE given the condition (4.18)

$$\sum_{i=1}^n e_i = 0. \quad (4.114)$$

- 5441 **4.6** Show that when $m = 1$, the confidence interval formula for the multiple linear
 5442 regression (4.102) is reduced to the confidence interval formula for the simple
 5443 linear regression (4.82).
- 5444 **4.7** Examine the global average December temperature anomalies from 1880 to
 5445 2018 in the dataset of the NOAAGlobalTemp.
 5446 (a) Make a linear regression of the temperature anomalies against time.
 5447 (b) Compute the confidence intervals of the fitted model at 95% confidence

- 5448 level.
- 5449 (c) Compute the confidence intervals of the anomaly data at 95% confidence
- 5450 level.
- 5451 (d) On the same figure similar to Fig. 4.6(a), plot the scatter plot of the
- 5452 anomaly data against time, and plot the confidence intervals computed in (b)
- 5453 and (c).
- 5454 **4.8** Make a diagnostic analysis for the regression results of the previous problem.
- 5455 (a) Produce a scatter plot of the residuals against time.
- 5456 (b) Visually check whether the assumptions of linearity and constant variance
- 5457 are satisfied.
- 5458 (c) Use the KS-test to check the normality assumption on residuals.
- 5459 (d) Use the DW-test to check the independence assumption on residuals.
- 5460 (e) When serial correlation is considered, find the effective degrees of freedom
- 5461 (edof).
- 5462 (f) Compute the confidence intervals in Steps (b) and (c) in the above
- 5463 problem using edof.
- 5464 (g) Produce a scatter plot of the anomaly data against time, and plot the
- 5465 confidence intervals on the same figure using the results of Step (f).
- 5466 **4.9** (a) Fit the NOAAGlobalTemp's global average annual mean data from 1880
- 5467 to 2018 to a ninth order orthogonal polynomial.
- 5468 (b) Plot the data and the fitted polynomial function on the same figure.
- 5469 (c) Produce a scatter plot of the residuals of the fitting against time.
- 5470 **4.10** Make a diagnostic analysis for the above regression and examine the regres-
- 5471 sion assumptions. In particular, use the KS-test to verify the normality as-
- 5472 sumption, and the DW-test to verify the independence assumption.
- 5473 **4.11** (a) Use multiple linear regression to compute the 12th order polynomial
- 5474 fitting of the NOAAGlobalTemp's global average annual mean data from 1880
- 5475 to 2018:
- $$T = b_0 + b_1 t + b_2 t^2 + \cdots + b_{12} t^{12} + \epsilon. \quad (4.115)$$
- 5476 (b) Plot the data and the fitted polynomial function on the same figure.
- 5477 (c) Produce a scatter plot of the residuals of the fitting against time.
- 5478 **4.12** Make a diagnostic analysis for the above regression and examine the regres-
- 5479 sion assumptions. In particular, use the KS-test to verify the normality as-
- 5480 sumption, and the DW-test to verify the independence assumption.
- 5481 **4.13** (a) Fit the global average **January** monthly mean temperature anomaly
- 5482 data from 1880 to 2018 in the NOAAGlobalTemp dataset to a third order
- 5483 orthogonal polynomial. The global average monthly mean NOAAGlobalTemp
- 5484 time series data are included in the book's master dataset named **data.zip**.
- 5485 You can also download the updated data from the Internet.
- 5486 (b) Plot the data and the fitted polynomial function on the same figure.
- 5487 (c) Produce a scatter plot of the residuals of the fitting against time in another
- 5488 figure.
- 5489 **4.14** Make a diagnostic analysis for the previous **January** data fit following the

- procedures in this chapter. In particular, use the KS-test to verify the normality assumption, and the DW-test to verify the independence assumption.
- 5490 procedures in this chapter. In particular, use the KS-test to verify the normality assumption, and the DW-test to verify the independence assumption.
- 5491
- 5492 **4.15** (a) Fit the global average **July** monthly mean temperature anomaly data from 1880 to 2018 in the NOAAGlobalTemp dataset to a third order orthogonal polynomial.
- 5493 (b) Plot the data and the fitted polynomial function on the same figure.
- 5494 (c) Produce a scatter plot of the residuals of the fitting against time in another figure.
- 5495
- 5496 **4.16** Make a diagnostic analysis for the previous **July** data fit following the pro-
- 5497 cedures in this chapter. In particular, use the KS-test to verify the normality
- 5500 assumption, and the DW-test to verify the independence assumption.
- 5501 **4.17** Use the gridded monthly NOAAGlobalTemp dataset and make a third-order
- 5502 polynomial fit to the **January** monthly mean temperature anomaly data from 1880 to 2018 in a grid box that covers Tokyo, Japan. The gridded
- 5503 monthly NOAAGlobalTemp dataset is included in the book's master data
- 5504 file `data.zip`. You can also download the updated data from the Internet.
- 5505
- 5506 **4.18** Use the gridded monthly NOAAGlobalTemp dataset and make a third-order
- 5507 polynomial fit to the **January** monthly mean temperature anomaly data from 1880 to 2018 in a grid box that covers Bonn, Germany.
- 5508
- 5509 **4.19** Use the gridded monthly NOAAGlobalTemp dataset and make a fifth-order
- 5510 polynomial fit to the monthly mean temperature anomaly data for a grid box,
- 5511 a month, and a period of time of your own interest. For example, you may
- 5512 choose your hometown grid box, the month you were born, and the period of
- 5513 1900-1999.
- 5514 **4.20** Make a diagnostic analysis for the previous data fit following the procedures
- 5515 in this chapter.
- 5516 **4.21** Use the January time series data of an USHCN station and fit a third or-
- 5517 der polynomial. Make a diagnostic analysis for the fit. The updated monthly
- 5518 USHCN station data may be downloaded from the Internet.

Appendix A A Tutorial of R and RStudio

5521

5522 The book uses R and the R Notebook. This chapter explains the installation of R
5523 and R Studio and demonstrates some basic uses of R.

5524 Equivalent Python codes and their Jupyter Notebooks may be found at the web-
5525 site www.climatemathematics.org.

5526 A.1 Download and install R and R-Studio

5527

5528 For Windows users, visit the website

5529 <https://cran.r-project.org/bin/windows/base/>

5530 to find the instructions for R program download and installations.

5531 For Mac users, visit

5532 <https://cran.r-project.org/bin/macosx/>

5533

5534 If you experience difficulties, please refer to online resources, Google or YouTube.
5535 A recent 3-minute YouTube instruction for R installation for Windows can be found
5536 from the following link:

5537 <https://www.youtube.com/watch?v=0hnk9hcx9M>

5538

5539 The same author also has a YouTube instruction about R installation for Mac
5540 (2 minutes):

5541 <https://www.youtube.com/watch?v=v=uxuuWXU-7UQ>

5542

5543 When R is installed, one can open R. The R Console window will appear. See
5544 Fig. A.1. One can use R Console to perform calculations, such as typing 2+3 and
5545 hitting return. However, most people today prefer using RStudio as the interface.

5546 To install RStudio, visit

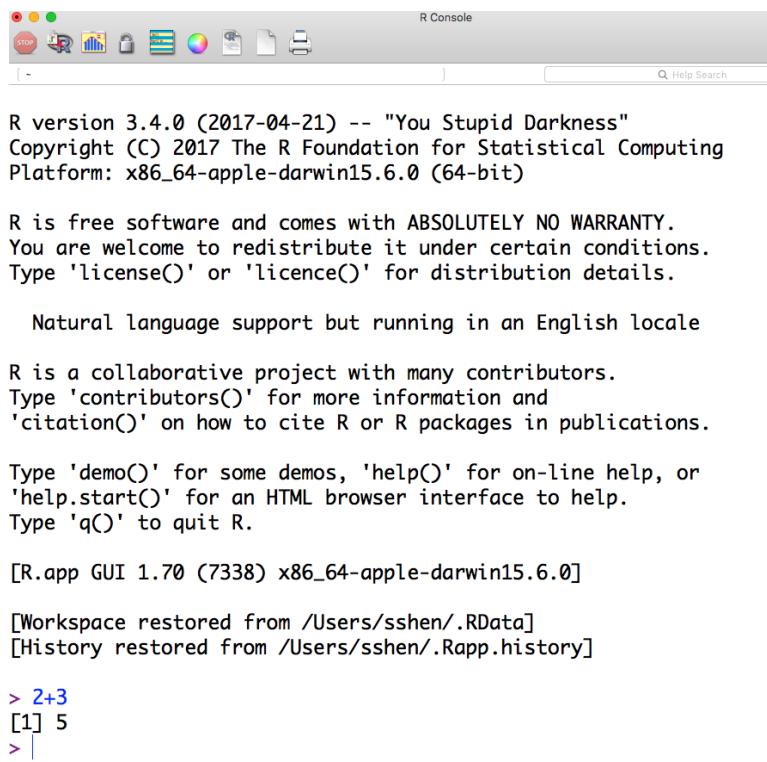
5547 <https://www.rstudio.com/products/rstudio/download/>

5548 This site allows one to choose Windows, or Mac OS, or Unix.

5549

5550 After both R and RStudio are installed, one can use either R or RStudio, or
5551 both, depending on one's interest. However, RStudio will not work without R.
5552 Thus, always install R first.

5553 When opening RStudio, four windows will appear as shown in Fig. A.2: The
5554 top left window is called R script, for writing the R code. The green arrow on



```
R version 3.4.0 (2017-04-21) -- "You Stupid Darkness"
Copyright (C) 2017 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[R.app GUI 1.70 (7338) x86_64-apple-darwin15.6.0]

[Workspace restored from /Users/sshen/.RData]
[History restored from /Users/sshen/.Rapp.history]

> 2+3
[1] 5
> |
```

Fig. A.1 R Console window after opening R.

5555 top of the window can be clicked to run the code. Each run is shown in the lower
 5556 left R Console window, and recorded on the upper right R History window. When
 5557 plotting, the figure will appear in the lower right R Plots window. For example,
 5558 `plot(x,x*x)` renders the eight points in the Plots window, because `x=1:8` defines
 5559 a sequence of numbers from 1 to 8. `x*x` yields a sequence from 1^2 to 8^2 .

A.2 R Tutorial

5560
 5561
 5562 Many excellent tutorials for quickly learning R programming, using a few hours or
 5563 a few evenings, are available online and in YouTube, such as
 5564 <http://ww2.coastal.edu/kingw/statistics/R-tutorials/>.
 5565 You can easily perform an Internet search and find your preferred tutorials.
 5566 It can be extremely difficult for the beginners of R to navigate through the offi-
 5567 cial, formal, detailed, and massive R-Project documentation:
 5568 <https://www.r-project.org/>
 5569

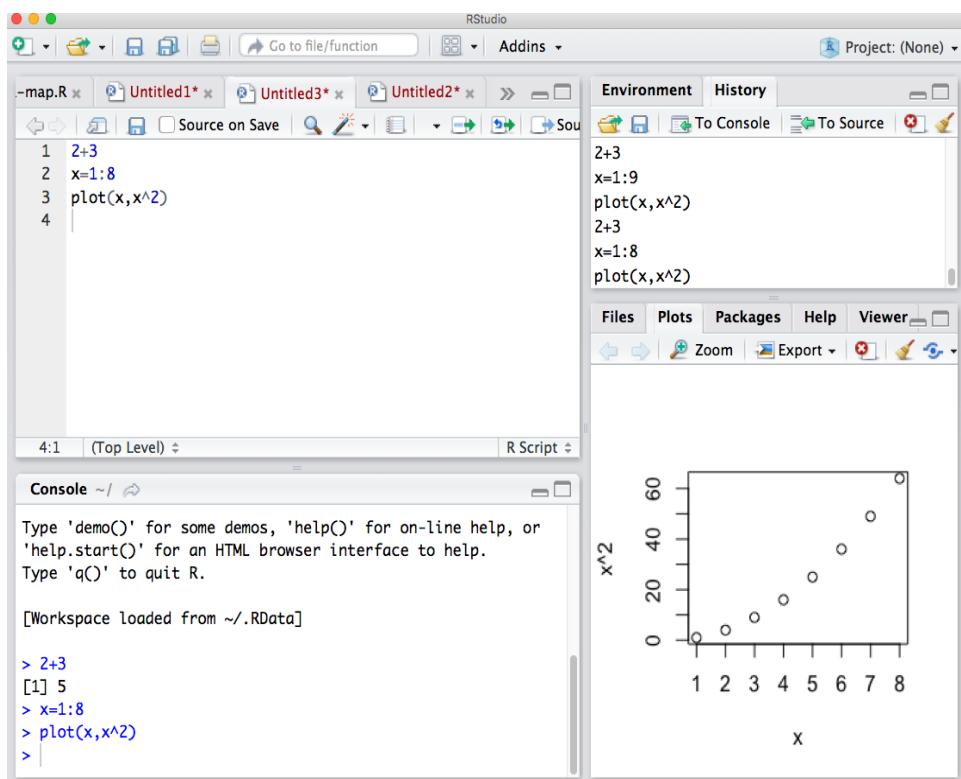


Fig. A.2 R Studio windows.

5570
5571

A.2.1 R as a smart calculator

5572 R can be used like a smart calculator that allows more elaborate calculations than
5573 those done with ordinary calculators.

```

5574 1+4
5575 #[1] 5 This is the result
5576 2+pi/4-0.8
5577 #[1] 1.985398
5578 x<-1
5579 y<-2
5580 z<-4
5581 t<-2*x^y-z
5582 t
5583 #[1] -2
5584 u=2      # "=" sign and "<-" are almost equivalent
5585 v=3      # The text behind the "#" sign is comments
5586 u+v
5587 #[1] 5

```

```
5588 sin(u*v)    # u*v = 6 in the sine function is considered radian by R
5589 #[1] -0.2794155
```

5590 R programming uses assignment operator `a <- b` to assign b to a. Often the equal
 5591 operator `a=b` can do the same job or vice versa. The two operators are equivalent
 5592 in general. However, certain R formulas have specific meanings for `=` and cannot be
 5593 replaced by `<-`. Most veteran R users use `<-` for assignment and `=` for defined R
 5594 formulas.

5595 A.2.2 Define a sequence in R

5597 Directly enter a sequence of daily maximum temperature data at San Diego International
 5598 Airport (Lat: 32.7336°N, Lon: 117.1831°W) during 1-7 May 2017 [unit:
 5599 °F].

5600 `tmax <- c(77, 72, 75, 73, 66, 64, 59)`
 5601 The data are from the Daily Summary of the Local Climatological Data (LCD),
 5602 National Centers for Environmental Information (NCEI)
 5603 <https://www.ncdc.noaa.gov/cdo-web/datatools/lcd>
 5604 The command `c()` is used to hold a data sequence and is named `tmax`. Entering
 5605 the `tmax` command will render the temperature data sequence:

```
5606 tmax
5607 #[1] 77 72 75 73 66 64 59
```

5608 You can generate different sequences using R, e.g.,
 5609 `1: 8 #Generates a sequence 1,2,...,8`

5610 Here the pound sign `#` begins R comments which are not executed by R calculations.
 5611 The same sequence can be generated by different commands, such as

```
5612 seq(1,8)
5613 seq(8)
5614 seq(1,8, by=1)
5615 seq(1,8, length=8)
5616 seq(1,8, length.out =8)
```

5617 The most useful sequence commands are `seq(1,8, by=1)` and `seq(1,8, length=8)`
 5618 or `seq(1,8, len=8)`. The former is determined by a begin value, end value, and
 5619 step size, and the latter by a begin value, end value, and number of values in the
 5620 sequence. For example, `seq(1951,2016, len=66*12)` renders a sequence of all the
 5621 months from January 1951 to December 2016.

5622
5623

A.2.3 Define a function in R

5624 The function command is
 5625 `name <- function(var1, var2, ...) expression of the function.`
 5626 For example,

```
5627 samfctn <- function(x) x*x
5628 samfctn(4)
5629 #[1] 16
5630 fctn2 <- function(x,y,z) x+y-z/2
5631 fctn2(1,2,3)
5632 #[1] 1.5
```

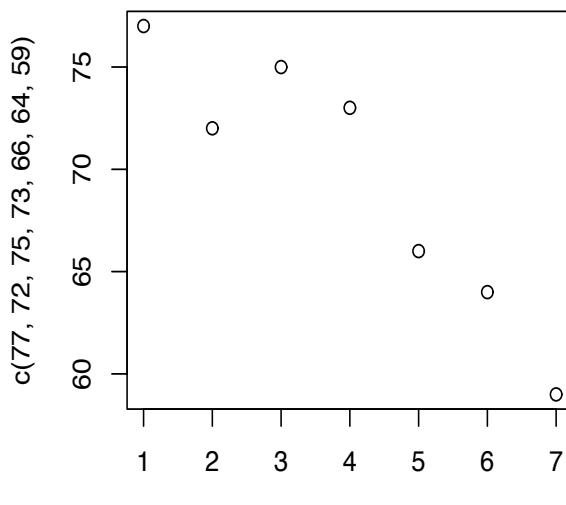
5633
5634

A.2.4 Plot with R

5635 R can plot all kinds of curves, surfaces, statistical plots, and maps. Below are a
 5636 few very simple examples for R beginners. For adding labels, ticks, color, and other
 5637 features to a plot, you will learn them from later parts of this book, and you may
 5638 also carry out an Internet search on R plot to find the commands for the proper
 5639 inclusion of the desired features.

5640 R plotting is based on the coordinate data. The following command plots the
 5641 seven days of San Diego Tmax data above:
 5642 `plot(1:7, c(77, 72, 75, 73, 66, 64, 59))`

The resulting graph is shown in Fig. A.3.

**Fig. A.3**

The daily maximum temperature during 1-7 May 2017 at San Diego International Airport.

5643

```

5644 plot(sin, -pi, 2*pi)    #plot the curve of y=sin(x) from -pi to 2 pi
5645
5646 square <- function(x) x*x    #Define a function
5647 plot(square, -3,2)    # Plot the defined function
5648
5649 # Plot a 3D surface
5650 x <- seq(-1, 1, length=100)
5651 y <- seq(-1, 1, length=100)
5652 z <- outer(x, y, function(x, y)(1-x^2-y^2))
5653 #outer (x,y, function) renders z function on the x, y grid
5654 persp(x,y,z, theta=330)
5655 # yields a 3D surface with perspective angle 330 deg
5656
5657 #Contour plot\index{contour plot}
5658 contour(x,y,z) #lined contours
5659 filled.contour(x,y,z) #color map of contours

```

5660 The color map of contours resulting from the last command is shown in Fig. A.4.

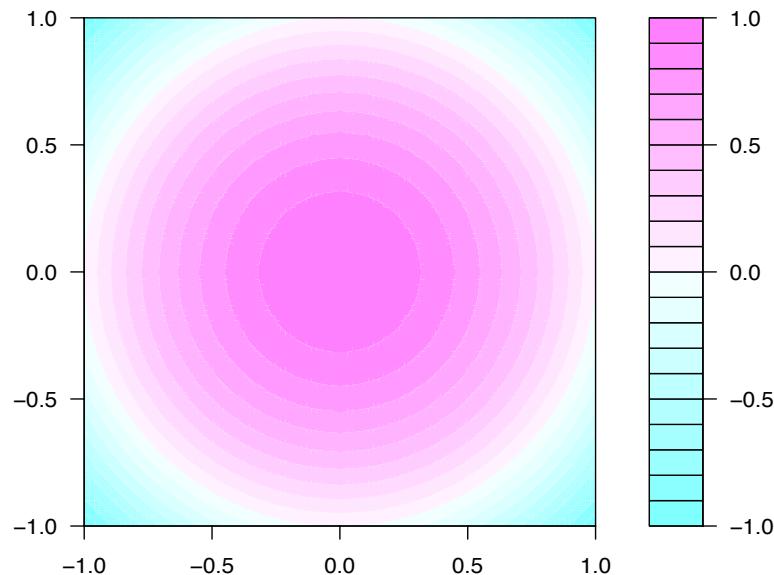


Fig. A.4 The color map of contours for the function $z = 1 - x^2 - y^2$.

5661

A.2.5 Symbolic calculations by R

5662 Many people once thought that R can handle only numbers. Actually R can also do
 5663 symbolic calculations, such as finding a derivative, although at present R is not the

5666 best symbolic calculation tool. One can use WolframAlpha, SymPy, and Yacas for
 5667 free symbolic calculations or use the paid software packages Maple or Mathematica.
 5668 Carry out an Internet search on symbolic calculation for calculus to find a long list
 5669 of symbolic calculation software packages, e.g.,
 5670 https://en.wikipedia.org/wiki/List_of_computer_algebra_systems.

```
5671 D(expression(x^2,'x'), 'x')
5672 # Take derivative of x^2 w.r.t. x
5673 2 * x #The answer is 2x
5674
5675 fx= expression(x^2,'x') #assign a function
5676 D(fx,'x') #differentiate the function w.r.t. x
5677 2 * x #The answer is 2x
5678
5679 fx= expression(x^2*sin(x),'x')
5680 #Change the expression and use the same derivative command
5681 D(fx,'x')
5682 2 * x * sin(x) + x^2 * cos(x)
5683
5684 fxy = expression(x^2+y^2, 'x','y')
5685 #One can define a function of 2 or more variables
5686 fxy #renders an expression of the function in terms of x and y
5687 #expression(x^2 + y^2, "x", "y")
5688 D(fxy,'x') #yields the partial derivative with respect to x: 2 * x
5689 D(fxy,'y') #yields the partial derivative with respect to y: 2 * y
5690
5691 square = function(x) x^2
5692 integrate (square, 0,1)
5693 #Integrate x^2 from 0 to 1 equals to 1/3 with details below
5694 #0.3333333 with absolute error < 3.7e-15
5695
5696 integrate(cos,0,pi/2)
5697 #Integrate cos(x) from 0 to pi/2 equals to 1 with details below
5698 #1 with absolute error < 1.1e-14
```

5699 The above two integration examples are for definite integrals. It seems that no
 5700 efficient R packages are available for finding antiderivatives, or indefinite integrals.

A.2.6 Vectors and matrices

5703 R can handle all kinds of operations involving vectors and matrices.

```
5704 c(1,6,3,pi,-3) #c() gives a vector and is considered a 4X1 column vector\index{column vect}
5705 #[1] 1.000000 6.000000 3.000000 3.141593 -3.000000
```

```
5706 seq(2,6) #Generate a sequence from 2 to 6
5707 #[1] 2 3 4 5 6
5708 seq(1,10,2) # Generate a sequence from 1 to 10 with 2 increment
5709 #[1] 1 3 5 7 9
5710 x=c(1,-1,1,-1)
5711 x+1 #1 is added to each element of x
5712 #[1] 2 0 2 0
5713 2*x #2 multiplies each element of x
5714 #[1] 2 -2 2 -2
5715 x/2 # Each element of x is divided by 2
5716 #[1] 0.5 -0.5 0.5 -0.5
5717 y=seq(1,4)
5718 x*y # This multiplication * multiplies each pair of elements
5719 #[1] 1 -2 3 -4
5720 x%*%y #This is the dot product of two vectors and yields
5721 # [,1]
5722 #[1,] -2
5723 t(x) # Transforms x into a row 1X4 vector
5724 # [,1] [,2] [,3] [,4]
5725 #[1,] 1 -1 1 -1
5726 t(x)%*%y #This is equivalent to dot product and forms 1X1 matrix
5727 # [,1]
5728 #[1,] -2
5729 x%*%t(y) #This column times row yields a 4X4 matrix\index{matrix}
5730 # [,1] [,2] [,3] [,4]
5731 #[1,] 1 2 3 4
5732 #[2,] -1 -2 -3 -4
5733 #[3,] 1 2 3 4
5734 #[4,] -1 -2 -3 -4
5735 my=matrix(y,ncol=2)
5736 #Convert a vector into a matrix of the same number of elements
5737 #The matrix elements\index{matrix!elements} go by column, the first column, second, etc
5738 #Command matrix(y,ncol=2, nrow=2), or matrix(y, ncol=2)
5739 #or matrix(y,2), or matrix(y,2,2) does the same job
5740 my
5741 # [,1] [,2]
5742 #[1,] 1 3
5743 #[2,] 2 4
5744 dim(my) #find dimensions of a matrix
5745 #[1] 2 2
5746
5747 bigM=matrix(1:100, nrow=10) #Generate a 10-by-10 matrix
5748 subM=bigM[4:6,3:7] #Extract a sub-matrix from a big matrix
5749 subM
```

```
5750 #      [,1] [,2] [,3] [,4] [,5]
5751 #[1,] 24   34   44   54   64
5752 #[2,] 25   35   45   55   65
5753 #[3,] 26   36   46   56   66
5754
5755 as.vector(my) #Convert a matrix to a vector, again via columns
5756 #[1] 1 2 3 4
5757 mx <- matrix(c(1,1,-1,-1), byrow=TRUE,nrow=2)
5758 mx*my #multiplication between each pair of elements
5759 #      [,1] [,2]
5760 #[1,] 1     3
5761 #[2,] -2    -4
5762 mx/my #division between each pair of elements
5763 #      [,1]      [,2]
5764 #[1,] 1.0  0.3333333
5765 #[2,] -0.5 -0.2500000
5766 mx-2*my
5767 #      [,1] [,2]
5768 #[1,] -1   -5
5769 #[2,] -5   -9
5770 mx%*%my #This is the real matrix multiplication\index{matrix multiplication} in matrix the
5771 #      [,1] [,2]
5772 #[1,] 3     7
5773 #[2,] -3   -7
5774 det(my) #determinant\index{determinant}
5775 #[1] -2
5776 myinv = solve(my) #yields the inverse of a matrix
5777 myinv
5778 #      [,1] [,2]
5779 #[1,] -2   1.5
5780 #[2,] 1    -0.5
5781 myinv%*%my #verifies the inverse of a matrix
5782 #      [,1] [,2]
5783 #[1,] 1     0
5784 #[2,] 0     1
5785 diag(my) #yields the diagonal vector of a matrix
5786 #[1] 1 4
5787 myeig=eigen(my) #yields eigenvalues\index{eigenvalue} and unit eigenvectors\index{eigenve
5788 myeig
5789 myeig$values
5790 #[1] 5.3722813 -0.3722813
5791 myeig$vectors
5792 #      [,1]      [,2]
5793 #[1,] -0.5657675 -0.9093767
```

```

5794 #[2,] -0.8245648 0.4159736
5795 mysvd = svd(my) #SVD decomposition of a matrix M=UDV'
5796 #SVD can be done for a rectangular matrix of mXn
5797 mysvd$d
5798 #[1] 5.4649857 0.3659662
5799 mysvd$u
5800 # [,1] [,2]
5801 #[1,] -0.5760484 -0.8174156
5802 #[2,] -0.8174156 0.5760484
5803 mysvd$v
5804 # [,1] [,2]
5805 #[1,] -0.4045536 0.9145143
5806 #[2,] -0.9145143 -0.4045536
5807
5808 ysol=solve(my,c(1,3))
5809 #solve linear equations\index{linear equations} matrix %*% x = b
5810 ysol #solve(matrix, b)
5811 #[1] 2.5 -0.5
5812 my%*%ysol #verifies the solution
5813 # [,1]
5814 #[1,] 1
5815 #[2,] 3

```

A.2.7 Simple statistics by R

5816
5817 R was originally designed to do statistical calculations. Thus, R has a comprehensive set of statistics functions and software packages. This sub-section gives a few basic commands. More will be described in the statistics chapter of this book.

```

5821 x=rnorm(10) #generate 10 normally distributed numbers
5822 x
5823 #[1] 2.8322260 -1.2187118 0.4690320 -0.2112469 0.1870511
5824 #[6] 0.2275427 -1.2619005 0.2855896 1.7492474 -0.1640900
5825 mean(x)
5826 #[1] 0.289474
5827 var(x)
5828 #[1] 1.531215
5829 sd(x)
5830 #[1] 1.237423
5831 median(x)
5832 #[1] 0.2072969
5833 quantile(x)
5834 # 0% 25% 50% 75% 100%
5835 #-1.2619005 -0.1994577 0.2072969 0.4231714 2.8322260

```

```

5836 range(x) #yields the min and max of x
5837 #[1] -1.261900 2.832226
5838 max(x)
5839 #[1] 2.832226
5840
5841 boxplot\index{boxplot}(x) #yields the box plot of x
5842 w=rnorm(1000)
5843
5844 summary(rnorm(12)) #statistical summary of the data sequence
5845 # Min. 1st Qu. Median Mean 3rd Qu. Max.
5846 #-1.9250 -0.6068 0.3366 0.2309 1.1840 2.5750
5847
5848 hist(w)
5849 #yields the histogram\index{histogram} of 1000 random numbers with a normal distribution
5850
5851 #Linear regression\index{linear regression} and linear trend\index{linear trend} line
5852 #2007-2016 data of the global temperature\index{anomalies!temperature} anomalies
5853 #Source: NOAA GlobalTemp data
5854 t=2007:2016
5855 T=c(.36,.30, .39, .46, .33, .38, .42, .50, .66, .70)
5856 lm(T ~ t) #Linear regression model of temp vs time
5857 #(Intercept) t
5858 #-73.42691 0.03673
5859 #Temperature change rate is 0.03673 deg C/yr or 0.37 deg C/decade
5860 plot(t,T, type="o", xlab="Year", ylab="Temperature [deg C]",
5861 main="2007-2016 Global Temperature\index{anomalies!temperature} Anomalies\index{anoma
5862 and Their Linear Trend [0.37 deg C/decade] ")
5863 abline(lm(T ~ t), lwd=2, col="red") #Regression line

```

5864 The global temperature data from 2007-2016 in the above R code example are
 5865 displayed in Fig. A.5, together with their linear trend line

$$T = -73.42691 + 0.03673t. \quad (\text{A.1})$$

5866

A.3 Online Tutorials

5867

5868 Numerous online R tutorials are available. Several are relatively efficient for learning
 5869 climate mathematics and are recommended below.

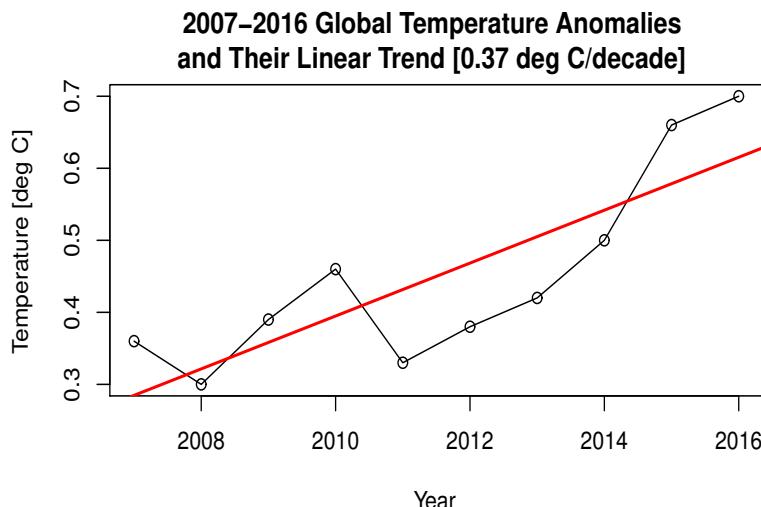


Fig. A.5 The 2007-2016 global average annual mean surface air temperature anomalies with respect to the 1971-2000 climate normal. The red line is a linear trend line computed from a linear regression model.

5870
5871

A.3.1 YouTube tutorial: for true beginners

5872
5873
5874

This is a very good and slow-paced 22-minute YouTube tutorial: Chapter 1. An Introduction to R
<https://www.youtube.com/watch?v=suVFuGET-OU>

5875
5876

A.3.2 YouTube tutorial: for some basic statistical summaries

5877
5878

This is a 9-minute tutorial by Layth Alwan.
<https://www.youtube.com/watch?v=Xj0ZQN-Nre4>

5879
5880

A.3.3 YouTube tutorial: Input data by reading a csv file into R

5881
5882
5883
5884

An excel file can be saved as csv file: xxxx.csv. This 15-minute YouTube video by Layth Alwan shows how to read a csv file into R. He also shows linear regression.
<https://www.youtube.com/watch?v=QkE8cp0B9gg>

R can input all kinds of data files, includingxlsx, txt, netCDF, MatLab data, Fortran file, and SAS data. Some commands are below. One can search the Internet to find proper data reading commands for any particular data format.

5885
5886
5887
5888
5889
5890

```
mydata <- read.csv("mydata.csv")
# read csv file named "mydata.csv"
```

```

5891 mydata <- read.table("mydata.txt")
5892 # read text file named "my data.txt"
5893
5894 library(gdata) # load the gdata package
5895 mydata = read.xls("mydata.xls") # read an excel file
5896
5897 library(foreign) # load the foreign package
5898 mydata = read.mtp("mydata.mtp") # read from .mtp file
5899
5900 library(foreign) # load the foreign package
5901 mydata = read.spss("myfile", to.data.frame=TRUE)
5902
5903 ff <- tempfile()
5904 cat(file = ff, "123456", "987654", sep = "\n")
5905 read.fortran(ff, c("F2.1","F2.0","I2")) #read a fotran file
5906
5907 library(ncdf4)
5908 ncin <- ncdf4::nc_open("ncfname") # open a NetCDF file
5909 lon <- ncvar_get(ncin, "lon") #read data "lon" from a netCDF file into R

```

5910 Many more details of reading and reformatting of .nc files will be discussed later
 5911 when dealing with NCEP/NCAR Reanalysis data.

5912 Some libraries are not in the R project. For example,

```

5913 library(ncdf4) #The following error message pops up
5914 Error in library(ncdf4) : there is no package called cdf4

```

5915 You can install the R package by

```

5916 install.packages("ncdf4")

```

5917 After this installation, `library(ncdf4)` will run, and the functions in the ncdf4
 5918 package will work.

5919 You only need to install the package once on your computer, but in each new R
 5920 session you must run `library(package)` in order to activate the package functions.
 5921 Many examples will be shown in the rest of this book.

5922 The R packages and the datasets used in this book are listed below and can be
 5923 downloaded and installed first before proceeding to the R codes in the rest of the
 5924 book.

```

5925 #R packages: animation, chron, e1071, fields, ggplot2, lattice,
5926 #latticeExtra, maps, mapdata, mapproj, matrixStats, ncdf4,
5927 #NLRoot, RColorBrewer, rgdal, rasterVis, raster, sp, TTR
5928
5929 #The zipped data:

```

```

5930 #https://cambridge.org/climate mathematics/data.zip
5931
5932 #To load a single package, such as "animation", you can do
5933 library(animation)
5934
5935 #You can also load all these packages in one shot using
5936 # pacman
5937
5938 install.packages("pacman")
5939 library(pacman)
5940 pacman::p_load(animation, chron, e1071, fields, ggplot2, lattice,
5941           latticeExtra, maps, mapdata, mapproj, matrixStats, ncdf4,
5942           NLRoot, RColorBrewer, rgdal, rasterVis, raster, sp, TTR)

5943 On your computer, you can create a directory called climmath under your user
5944 name. The one used in the book is Users/sshens/climmath. You unzip the data and
5945 move the data folder under the Users/sshens/climmath directory. A data folder will
5946 created: Users/sshens/climmath/data. The data folder contains about 400 MB of
5947 data. Place all the R codes in the directory Users/sshens/climmath. Then, you can
5948 run all the codes in this book after replacing sshens by your user name on your own
5949 computer.

```

A.4 Chapter summary

5950 This chapter has described the following basic aspects of the R programming language
 5951
 5952 (i) Installation of both R and R Studio.
 5953 (ii) Layout of R Studio panels and functions.
 5954 (iii) R commands for defining and computing vectors, matrices, and functions.
 5955 (iv) Simple statistics using R: mean, median, standard deviation, variance, histogram, boxplot, scatterplot, and linear regression.
 5956 (v) Online resources for R tutorial material.
 5957
 5958 We suggest that reading and practicing the R basics included this chapter might
 5959 require about three hours. You might then use perhaps five to ten hours to do
 5960 the exercise problems of this chapter. After that introduction to R, we estimate
 5961 that you should be able to develop and carry out projects using R independently.
 5962 Familiarizing yourself with additional R commands may require spending some time
 5963 online with a search engine.
 5964 More sophisticated examples of using R for the analysis and visualization of the
 5965 data from climate models and observations are included in Chapters 9-11. These
 5966
 5967

5968 chapters also provide research-level R-graphics codes for climate sciences, R pro-
5969 grams for advanced statistical analysis of climate data, such as empirical orthogonal
5970 functions computed from the climate model datasets, and R programming tech-
5971 niques for analyzing datasets with missing data, either in space or time. To develop
5972 R projects at a more advanced level, you will need to read these chapters or search
5973 for the R code for the specific analysis or graphics tasks described in these chapters.
5974 These tasks include examples such as the preparation of data from a global climate
5975 model (GCM) for a singular value decomposition (SVD) analysis.

5976 The R programming language was created by statisticians Ross Ihaka and Robert
5977 Gentleman in New Zealand and was first released in 1995. R is named partly after
5978 the first names of the two original two authors of R, and partly as a play on the
5979 name of the S programming language for statistics. Currently, R is a popular com-
5980 puter programming language used in almost every field of science and engineering.
5981 Among the top programming languages ranked by the Institute of Electric and
5982 Electronic Engineers (IEEE) in 2017, R ranks sixth, following Python, C, Java,
5983 C++, and C#. For a climate science student or professional who is not a specialist
5984 in computer programming or information technology, R is easy to learn, and it offers
5985 numerous public-domain software packages that are free to use. As an alternative
5986 to R, equivalent Python codes for our entire book are available online at the book
5987 website www.climatemathematics.org.

References and Further Readings

5989 [1] King, K.B., 2016: *R Tutorials*, Coastal Carolina University. URL:

5990 <http://ww2.coastal.edu/kingw/statistics/R-tutorials/>

This easy-to-learn tutorial is for beginners of R, and does not require any computer programming background. It contains many statistics examples.

5991

5992 [2] Jost, S., 2018: *Introduction to R: An R Tutorial for Data Analysis and Regression*. De Paul University. URL:

5993 <http://facweb.cs.depaul.edu/sjost/csc423/>

This is a very brief R tutorial from the perspective of computer programming for beginners. One needs only about one hour to go through the entire tutorial.

5995

5996 [3] Vallis, G.K., 2016: Geophysical fluid dynamics: whence, whither
5997 and why? *Proceedings of the Royal Society A*, 472: 20160140.
5998 <http://dx.doi.org/10.1098/rspa.2016.0140>

In this stimulating article, Vallis discusses the role of geophysical fluid dynamics in understanding the dynamics of atmospheres and oceans. Geoffrey K. Vallis, a Professor in the Department of Mathematics at the University of Exeter, is an expert in climate dynamics, the circulation of planetary atmospheres, and dynamical meteorology and oceanography. He is the author of *Atmospheric and Oceanic Fluid Dynamics: Fundamentals and Large-Scale Circulation*. This widely praised standard text is a magisterial treatment of geophysical fluid dynamics. The book was first published in 2006, and its second edition in 2017 contains nearly 1,000 pages.

5999

6000

Exercises

6001

6002 **A.1** Use R to define a data sequence `t=seq(2015,2018, length=100)`, and then
 6003 plot the following two functions on the same figure: $y = \sin(2\pi(t - 0.1))$ and
 6004 $y = \cos^2(2\pi t)$.

6005 **A.2** (a) Use R to make a contour plot of the function $z = \sin^2 x \cos^2(y - \pi)$ over
 6006 the domain of $[0, 2\pi] \times [0, 2\pi]$.
 6007 (b) Use R to plot a color contour map for the same function on the same
 6008 domain.

6009 **A.3** Use R to solve the following linear equations:

$$\begin{cases} 9x + 8y = 87 \\ 6x - 20y = 126 \end{cases}$$

6010 **A.4** Use R to solve the following linear equations:

$$\begin{cases} -3x + 2y + z = 1 \\ -2x - y + z = 2 \\ 2x + y - 4z = 0 \end{cases}$$

6011 **A.5** For some purposes, climatology or climate is defined as the mean state, or normal state, of a climate parameter, and is calculated from data over a period of time called the climatology period (e.g., 1961-1990). Thus the surface air temperature climate or climatology at a given location may be calculated by averaging observational temperature data over a period such as 1961 through 1990. Thirty years are often considered in the climate science community as the standard length of a climatology period. Due to the relatively high density of weather stations in 1961-1990, compared to earlier periods, investigators have often used 1961-1990 as their climatology period, although some may now choose 1971-2000 or 1981-2010. Surface air temperature (SAT) is often defined as the temperature inside a white-painted louvered instrument container or box, known as a Stevenson screen located on a stand about 2 meters above the ground. The purpose of the Stevenson screen is to shelter the instruments from radiation, precipitation, animals, leaves, etc, while allowing the air to circulate freely inside the box. The daily maximum temperature (Tmax) is the maximum temperature measured inside the screen box by a maximum temperature thermometer within 24 hours. The daily minimum temperature (Tmin) is the minimum temperature within 24 hours. The daily mean temperature (Tmean) is the average of Tmax and Tmin.

6030 Go to the USHCN website

6031 http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn_map_interface.html
 6032

6033 and download the monthly Tmax, Tmin, and Tmean data of the Cuyamaca

6034 station (USHCN Site No. 042239) near San Diego, California, USA. Or down-
6035 load the data of this book at

6036 www.cambridge.org/climate-mathematics/data.zip
6037

6038 Unzip the data and find the data from the file named **CA042239T.csv**

6039 (a) Use R to arrange the monthly Cuyamaca Tmax data from January 1961
6040 to December 1990 as a matrix with each row as year and each column as
6041 month.

6042 (b) Do the same for Tmin.
6043 (c) Do the same for Tmean.

6044 **A.6** (a) Use R to calculate the August climatology of Tmax, Tmin, and Tmean
6045 for the Cuyamaca station according to the 1961-1990 climatology period.

6046 (b) Use R to compute the standard deviation of Tmax, Tmin, and Tmean of
6047 the Cuyamaca station for January during the 1961-1990 climatology period.

6048 **A.7** (a) Use R to plot the the Cuyamaca January Tmin time series from 1951 to
6049 2010 with a continuous curve.

6050 (b) Use R to plot the linear trend lines of Tmin on the same plot as (a) in
6051 the following time periods:

6052 (i) 1951-2010,
6053 (ii) 1961-2010,
6054 (iii) 1971-2010, and
6055 (iv) 1981-2010.

6056 (c) Finally, what is the temporal trend per decade for each of the four periods
6057 above?

6058 **A.8** Use R to plot the time series and its trend line for P.D. Jones' global average
6059 annual mean temperature anomaly data: **JonesGlobalT.txt**. This data file
6060 can be found from the book's data.zip file downloaded from the book website.

6061 (a) Plot the global average annual mean temperature from 1880 to 2015.
6062 (b) Find the linear trend of the temperature from 1880 to 2015. Plot the
6063 trend line on the same figure as (a).

6064 (c) Find the linear trend from 1900 to 1999. Plot the trend line on the same
6065 figure as (a).

6066 **A.9** Use the gridded NOAA global monthly temperature anomaly data NOAA-
6067 GlobalTemp from the following website or another data source

6068 [https://www.ncdc.noaa.gov/data-access/marineocean-data/
6069 noaa-global-surface-temperature-noaaglobaltemp](https://www.ncdc.noaa.gov/data-access/marineocean-data/noaa-global-surface-temperature-noaaglobaltemp)

6070 Or use the **NOAAGlobalT.csv** data file from the book's data.zip file down-
6071 loaded from the book website. Choose two 5-by-5 degrees lat-lon grid boxes
6072 of your interest. Plot the temperature anomaly time series of the two boxes
6073 on the same figure using two different colors.

- 6074 **A.10** Use the same NOAAGlobalTemp dataset, choose sufficiently many grid boxes
6075 that cover the state of Texas, USA. Compute the average temperature anomalies
6076 of these boxes for each month. Then plot the monthly average temperature
6077 anomalies as a function of time. Plot a linear trend line on the same figure.
- 6078 **A.11** Choose a 5-by-5 degrees grid box that covers Edmonton, Canada, and another
6079 grid box that covers San Diego, USA.
6080 (a) Use R and 30 years of the January NOAAGlobalTemp data from January
6081 1981 to January 2010 to compute the standard deviations for each grid box
6082 for January.
6083 (b) Do the same for February, March, ..., December.
6084 (c) Use R to write your standard deviation results in a 12-by-2 matrix with
6085 each row for a month, and each column for a grid box ID.
6086 (d) Describe the main differences between the values of the two columns.

B Appendix B Visualization of Matrices

6089

6088 People talk about climate data frequently, also read or imagine climate data, and
 6091 yet rarely play with them and use them, because people often think that it takes
 6092 a computer expert to do that. However, that has changed. With today's technol-
 6093 ogy, now anyone can use a computer to play with climate data, such as a sequence
 6094 of temperature values of a weather station at different observed times, a matrix
 6095 of data for a station for temperature, air pressure, precipitation; wind speed, and
 6096 wind direction at different times; and an array of temperature data on a 5-degree
 6097 latitude-longitude grid for the entire world for different months. The first is a vec-
 6098 tor. The second is a variable-time matrix, and a space-time 3-dimensional array.
 6099 When considering temperature variation in time at different air pressure levels and
 6100 different water depth, we need to add one more dimension: the altitude. The tem-
 6101 perature data for ocean and atmosphere for the Earth is a 4-dimensional array,
 6102 with 3D space and 1D time. This chapter attempts to provide basic statistical and
 6103 computing methods to describe and visualize some simple climate datasets. As the
 6104 book progresses, more complex statistics and data visualization will be introduced.

6105 We use both R and Python computer codes in this book for computing and
 6106 visualization. Our method description is stated in R. A Python code following each
 6107 R code is included in a box with a light yellow background. You can also learn
 6108 the two computer languages and their applications to climate data from the book
 6109 "*Climate Mathematics: Theory and Applications*" (Shen and Somerville 2019) and
 6110 its website www.climatemathematics.org.

6111 The climate data used in this book are included in the `data.zip` file download-
 6112 able from our book website www.climatestatistics.org. You can also obtain the
 6113 updated data from the original data providers, such as www.esrl.noaa.gov and
 6114 www.ncei.noaa.gov.

6115 After learning this chapter, a reader should be able to analyze simple climate
 6116 datasets, compute data statistics, and plot the data in various ways.

6117

6118

6119

B.1 Global temperature anomalies from 1880 to 2018: data visualization and statistical indices

6120

In a list of popular climate datasets, the global average annual mean surface air
 6121 temperature anomalies might be on top. Here, the *anomalies* means the temperature

6122 departures from the normal temperature that is called *climatology*. Climatology is
 6123 usually defined as the mean of temperature data in a given period of time, such as
 6124 from 1971 to 2020. Thus, the temperature anomaly data are the differences of the
 6125 temperature data minus the climatology.

6126 This section will use the global average annual mean surface air temperature
 6127 anomaly dataset as an example to describe some basic statistical and computing
 6128 methods.

6129 B.1.1 The NOAA GlobalTemp dataset

6131 The 1880-2018 global average annual mean surface air temperature (SAT) anomaly
 6132 data are shown as follows:

```
6133 [1] -0.37 -0.32 -0.32 -0.40 -0.46 -0.47 -0.45 -0.50 -0.40 -0.35 -0.57
6134 [12] -0.49 -0.55 -0.56 -0.53 -0.47 -0.34 -0.36 -0.50 -0.37 -0.31 -0.39
6135 [23] -0.49 -0.58 -0.66 -0.53 -0.46 -0.62 -0.69 -0.68 -0.63 -0.68 -0.58
6136 [34] -0.57 -0.39 -0.32 -0.54 -0.56 -0.45 -0.45 -0.46 -0.39 -0.47 -0.46
6137 [45] -0.50 -0.39 -0.31 -0.40 -0.42 -0.54 -0.34 -0.32 -0.36 -0.49 -0.35
6138 [56] -0.38 -0.36 -0.26 -0.27 -0.26 -0.15 -0.05 -0.09 -0.09 0.04 -0.08
6139 [67] -0.25 -0.30 -0.30 -0.30 -0.41 -0.26 -0.22 -0.15 -0.36 -0.38 -0.44
6140 [78] -0.19 -0.13 -0.18 -0.22 -0.16 -0.15 -0.13 -0.39 -0.32 -0.27 -0.26
6141 [89] -0.27 -0.15 -0.21 -0.32 -0.22 -0.08 -0.32 -0.24 -0.32 -0.05 -0.13
6142 [100] -0.02 0.02 0.06 -0.06 0.10 -0.10 -0.11 -0.01 0.13 0.13 0.05
6143 [111] 0.19 0.16 0.01 0.03 0.09 0.21 0.08 0.27 0.39 0.20 0.18
6144 [122] 0.30 0.35 0.36 0.33 0.41 0.37 0.36 0.30 0.39 0.45 0.33
6145 [133] 0.38 0.42 0.50 0.66 0.70 0.60 0.54
```

6146 The data are part of the dataset named as the NOAA Merged Land Ocean Global
 6147 Surface Temperature Analysis (NOAA GlobalTemp) V4. The dataset was generated
 6148 by the NOAA National Centers for Environmental Information (NCEI) (Smith et
 6149 al. 2008; Vose et al. 2012). Here, NOAA stands for the United States National
 6150 Oceanic and Atmospheric Administration.

6151 The anomalies are with respect to the 1971-2000 climatology, i.e., 1971-2000
 6152 mean. An anomaly of a weather station datum is defined by the datum minus the
 6153 station climatology. The first anomaly datum -0.37°C , indexed by [1], in the
 6154 above data table corresponds to 1880 and the last to 2018, a total of 139 years. The
 6155 last row is indexed from [133] to [139].

6156 One might be interested in various kinds of statistical characteristics of the data,
 6157 such as mean, variance, standard deviation, skewness, kurtosis, median, 5th per-
 6158 centile, 95th percentile, and other quantiles. Is the data's probabilistic distribution
 6159 approximately normal? What does the box plot look like? Are there any outliers?
 6160 What is a good graphic representation of the data, i.e., popularly known as a climate
 6161 figure?

6162 When considering global climate changes, why do scientists often use anomalies,

6163 instead of the full values directly from the observed thermometer readings? This is
 6164 because that the observational estimates of the global average annual mean surface
 6165 temperature are less accurate than the similar estimates for the changes from year
 6166 to year. There is a concept of characteristic spatial correlation length scale for a
 6167 climate variable, such as surface temperature. The length scale is often computed
 6168 from anomalies.

6169 The use of anomalies is also a way of reducing or eliminating individual station
 6170 biases. A simple example of such biases is that due to station location, which is
 6171 usually fixed in a long period of time. It is easy to understand, for instance, that
 6172 a station located in the valley of a mountainous region might report surface tem-
 6173 peratures that are higher than the true average surface temperature for the entire
 6174 region and cannot be used to describe the behavior of climate change in the region.
 6175 However, the anomalies at the valley station may synchronously reflect the charac-
 6176 teristics of the anomalies for the region. Many online materials give justifications
 6177 and examples on the use of climate data anomalies, e.g., NOAA NCEI (2021).

6178 The global average annual mean temperature anomalies quoted above are also
 6179 important for analyzing climate simulations. When we average over a large scale,
 6180 many random errors cancel out. When we investigate the response of such large
 6181 scale perturbations as the variations of Sun's brightness or atmospheric carbon
 6182 dioxide, these averaged data can help validate and improve climate models. See the
 6183 examples in the book by Hennemuth et al. (2013) that includes many statistical
 6184 analyses of both observed and model data.

6185 **B.1.2 Visualize the data of global average annual mean** 6186 **temperature**

6188 Many different ways have been employed to visualize the global average annual
 6189 mean temperature anomalies. The following three are popular ones appearing in
 6190 scientific and news publications: (a) a simple point-line graph, (b) a curve of stair-
 6191 case steps, and (c) a color bar chart, as shown in Figs. B.1 - B.3. This subsection
 6192 shows how to generate these figures by R and Python computer programming lan-
 6193 guages. The Python codes are in yellow boxes.

6194 **B.1.2.1 Plot a point-line graph of time series data**

6195 Figure B.1 is a simple line graph that connects all the data points, denoted by
 6196 the small circles, together to form a curve showing the historical record of the
 6197 global temperature anomalies. It is plotted from the the NOAAGlobalTemp V4
 6198 data quoted above. The figure can be generated by the following R code.

```
6199 # R plot Fig. 1.1: A simple line graph of data
6200 # go to your working directory
6201 setwd("/Users/sshen/climstats")
6202 # read the data file from the folder named "data"
6203 NOAAtemp = read.table(
  "data/aravg.ann.land_ocean.90S.90N.v4.0.1.201907.txt",
```

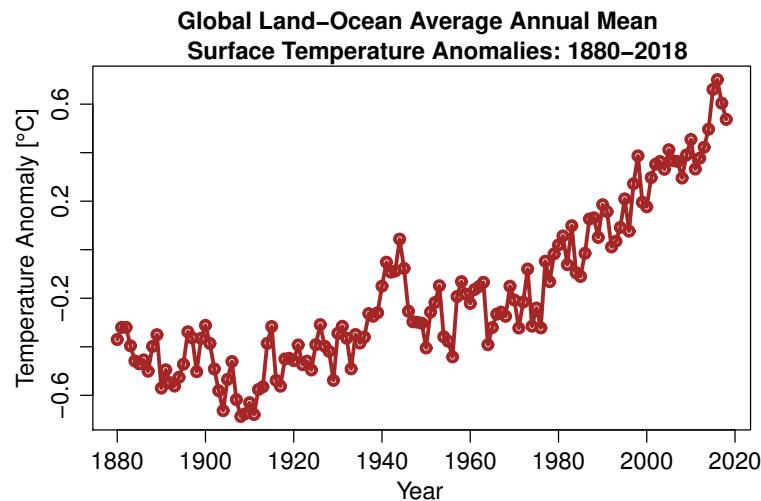


Fig. B.1 Point-line graph of the 1880–2018 global average annual mean temperature anomalies with respect to the 1971–2000 climatology, based on the NOAAGlobalTemp V4 data.

```

6205     header=FALSE) #Read from the data folder
6206     dim(NOAAtemp) # check the data matrix dimension
6207     #[1] 140   6 #140 years from 1880 to 2019
6208     #2019 will be excluded since data only up to July 2019
6209     #col1 is year, col2 is anomalies, col3-6 are data errors
6210     par(mar=c(3.5,3.5,2.5,1), mgp=c(2,0.8,0))
6211     plot(NOAAtemp[1:139,1], NOAAtemp[1:139,2],
6212       type = "l", col="brown", lwd=3, cex.lab=1.2, cex.axis=1.2,
6213       main ="Global_Land_Ocean_Average_Annual_Mean
6214       Surface_Temperature_Anomalies_1880-2018",
6215       xlab="Year",
6216       ylab=expression(
6217         paste("Temperature_Anomaly[", degree, "C]")))

```

```

# Python plot Fig. 1.1: A simple line graph of data
# Go to your working directory
os.chdir("/Users/sshenn/climstats")
# read the data file from the folder named "data"
NOAAtemp = read_table(\n
    "data/aravg.ann.land_ocean.90S.90N.v4.0.1.201907.txt",
    header = None, delimiter = "\s+")
# check the data matrix dimension
print("The dimension of our data table is:", NOAAtemp.shape)
x = np.array(NOAAtemp.loc[0:138, 0])
y = np.array(NOAAtemp.loc[0:138, 1])
plt.plot(x, y, 'brown', linewidth = 3);
plt.title("Global Land-Ocean Average Annual Mean\nSurface Temperature Anomaly: 1880-2018", pad = 15)
plt.xlabel("Year", size = 25, labelpad = 20)
plt.ylabel("Temperature Anomaly [$\degree C]", size = 25, labelpad = 20)
plt.show() # display on screen

```

6218

6219

B.1.2.2 Plot a staircase chart

6220 The staircase chart Fig. B.2 shows both data and the year-to-year annual changes
 6221 of temperature. One can clearly see that some years have a larger change from
 6222 their previous year, while other years have a smaller change. This information can
 6223 be used for further climate analysis.

6224 Denote the global average annual mean temperature by $T(t)$ where t stands for
 6225 time in year, the 1971-2000 climatology is by C , and the anomalies by $A(t)$. We
 6226 have

$$T(t) = C + A(t). \quad (\text{B.1})$$

6227 The temperature change from its previous year is

$$\Delta T = T(t) - T(t-1) = A(t) - A(t-1) = \Delta A. \quad (\text{B.2})$$

6228 This implies that an anomaly change is equal to its corresponding temperature
 6229 change. We do not need to evaluate the 1971-2000 climatology C of the global
 6230 average annual mean, when studying changes, as discussed earlier.

6231 Figure B.2 can be generated by the following R code

```

#R plot Fig. 1.2: Staircase chart of data
plot(NOAAtemp[1:139,1], NOAAtemp[1:139,2],
     type="s", #staircase curve for data
     col="black", lwd=2,
     main="Global Land-Ocean Average Annual Mean
           Surface Temperature Anomalies: 1880-2018",
     cex.lab=1.2,cex.axis=1.2,
     xlab="year",
     ylab=expression(paste(
           "Temperature Anomaly [", degree, "C]")))

```

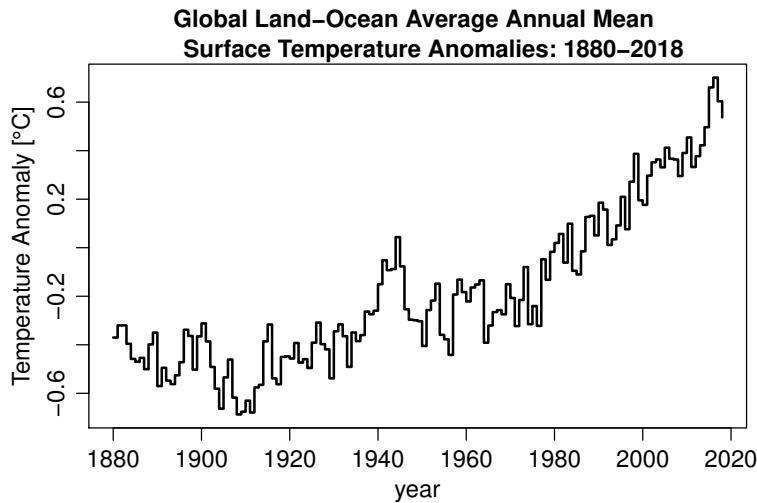


Fig. B.2 Staircase chart of the 1880–2018 global average annual mean temperature anomalies.

6243 The corresponding Python code is in the following box.

```
# Python plot Fig. 1.2: A staircase chart of data
# keyword arguments
kwargs = {'drawstyle' : 'steps'}
plt.plot(x, y, 'black', linewidth = 3, **kwargs);
plt.title("Global_Land-Ocean_Average_Annual_Mean\nSurface_Temperature_Anomaly: 1880 - 2018", pad = 15)
plt.xlabel("Year", size = 25, labelpad = 20)
plt.ylabel("Temperature_Anomaly [°C]", size = 25, labelpad = 20)
plt.show()
```

6244

6245 B.1.2.3 Plot a bar chart with colors

6246 Figure B.3 shows the anomaly size for each year by a color bar: red for positive
 6247 anomalies and blue for negative anomalies. This bar chart style has an advantage of
 6248 visualizing climate extremes, since these extremes may leave a distinct impression
 6249 on viewers. The black curve is a 5-point moving average of the annual anomaly data,
 6250 computed for each year by the mean of that year, together with two years before
 6251 and two years after. The moving average smooths the high frequency fluctuations
 6252 to reveal the long-term trends of temperature variations.

```
# R plot Fig. 1.3: A color bar chart of data
x <- NOAAtemp[,1]
y <- NOAAtemp[,2]
z <- rep(-99, length(x))
# compute 5-point moving average
```

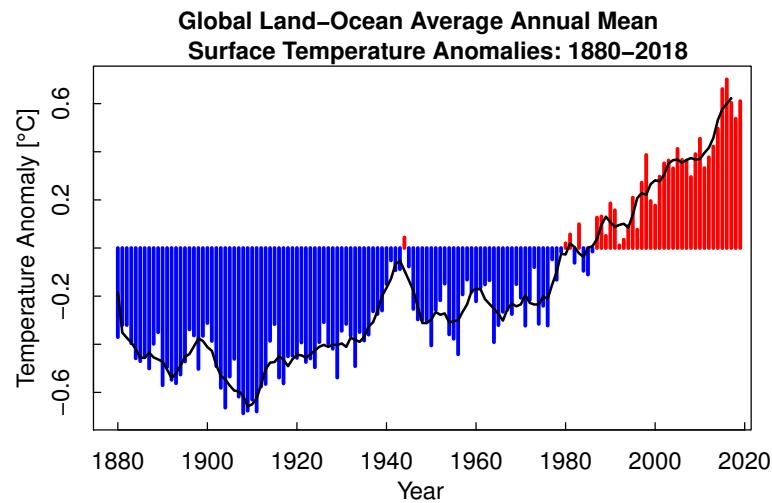


Fig. B.3 Color bar chart of the 1880-2018 global average annual mean temperature anomalies.

```

6258 for (i in 3:length(x)-2) z[i] =
6259   mean(c(y[i-2],y[i-1],y[i],y[i+1],y[i+2]))
6260 n1 <- which(y>=0); x1 <- x[n1]; y1 <- y[n1]
6261 n2 <- which(y<0); x2 <- x[n2]; y2 <- y[n2]
6262 x3 <- x[2:length(x)-2]
6263 y3 <- z[2:length(x)-2]
6264 plot(x1, y1, type="h", #bars for data
6265   xlim = c(1880,2016), lwd=3,
6266   tck = 0.02, #tck>0 makes ticks inside the plot
6267   ylim = c(-0.7,0.7), xlab="Year", col="red",
6268   ylab = expression(paste(
6269     "Temperature", degree, "C")))
6270 main ="Global", "Land", "-", "Ocean", "Average", "Annual", "Mean"
6271 "Surface", "Temperature", "Anomalies:", "1880", "-", "2018",
6272 cex.lab = 1.2, cex.axis = 1.2)
6273 lines(x2, y2, type="h",
6274   lwd = 3, tck = -0.02, col = "blue")
6275 lines(x3, y3, lwd = 2)

```

```

# Python plot Fig. 1.3: A color bar chart of data
# define an array z with y number of ones
z = np.ones(y.size)
z[0] = -99
z[1] = -99
# computer 5-point moving average
for i in range(2, z.size - 2):
    z[i] = np.mean(y[i-2:i+3])
z[z.size - 2] = -99
z[z.size - 1] = -99
# define variables based on range
y1 = [y[i] for i in range(y.size) if y[i] >= 0]
x1 = [x[i] for i in range(y.size) if y[i] >= 0]
y2 = [y[i] for i in range(y.size) if y[i] < 0]
x2 = [x[i] for i in range(y.size) if y[i] < 0]
y3 = z[2:x.size-2]
x3 = x[2:x.size-2]

# plot the moving average
plt.plot(x3, y3, 'black', linewidth = 3)
plt.title("Global_Land-Ocean_Average_Annual_Mean\nSurface_Temperature_Anomaly_1880-2018", pad = 20)

# create bar chart
plt.bar(x1, y1, color = 'red');
plt.bar(x2, y2, color = 'blue');
plt.xlabel("Year", size = 25, labelpad = 20)
plt.ylabel("Temperature_Anomaly[$\degree$C]", size = 25, labelpad = 20)
plt.show()

```

6276

6277
6278

B.1.3 Statistical indices

6279 The commonly used basic statistical indices include mean, variance , standard de-
 6280 viation, skewness, kurtosis, and quantiles. We first use R to calculate these indices
 6281 for the global average annual mean temperature anomalies. Then we describe their
 6282 mathematical formulas and interpret the numerical results.

```

#R code for computing statistical indices
setwd("/Users/sshen/climstats")
NOAAtemp = read.table(
  "data/aravg.ann.land_ocean.90S.90N.v4.0.1.201907.txt",
  header=FALSE)
temp2018=NOAAtemp[1:139,2] #use the temp data up to 2018
head(temp2018) #show the first six values
#[1] -0.370221 -0.319993 -0.320088 -0.396044 -0.458355 -0.470374
mean(temp2018) #mean
#[1] -0.1858632
sd(temp2018) #standard deviation
#[1] 0.324757
var(temp2018) #variance
#[1] 0.1054671

```

```

6297 library(e1071)
6298 #This R library is needed to compute the following parameters
6299 #install.packages("e1071") #if it is not in your computer
6300 skewness(temp2018)
6301 #[1] 0.7742704
6302 kurtosis(temp2018)
6303 #[1] -0.2619131
6304 median(temp2018)
6305 #[1] -0.274434
6306 quantile(temp2018,probs= c(0.05, 0.25, 0.75, 0.95))
6307 #      5%       25%       75%       95%
6308 #   -0.5764861 -0.4119770  0.0155245  0.4132383

```

6309 We use $x = \{x_1, x_2, \dots, x_n\}$ to denote the sampling data for a time series. The
 6310 statistical indices computed above by R are based on the following mathematical
 6311 formulas for mean, variance, standard deviation, skewness, and kurtosis:

$$\text{Mean: } \mu(x) = \frac{1}{n} \sum_{k=1}^n x_k, \quad (\text{B.3})$$

$$\text{Variance by unbiased estimate: } \sigma^2(x) = \frac{1}{n-1} \sum_{k=1}^n (x_k - \mu(x))^2, \quad (\text{B.4})$$

$$\text{Standard deviation: } \sigma(x) = (\sigma^2(x))^{1/2}, \quad (\text{B.5})$$

$$\text{Skewness: } \gamma_3(x) = \frac{1}{n} \sum_{k=1}^n \left(\frac{x_k - \mu(x)}{\sigma} \right)^3, \quad (\text{B.6})$$

$$\text{Kurtosis: } \gamma_4(x) = \frac{1}{n} \sum_{k=1}^n \left(\frac{x_k - \mu(x)}{\sigma} \right)^4 - 3. \quad (\text{B.7})$$

6312 A *quantile* cuts a sorted sequence of data. For example, the 25th quantile, also
 6313 called 25th percentile or 25% quantile, is the value that has 25% of the sorted data
 6314 smaller than this value and 75% larger than this value. The 50th percentile is also
 6315 known as the median.

```

# Python code for computing statistical indices
temp2018 = np.array(NOAAtemp.loc[0:138, 1]) # data string
# arithmetic average of the data
print("The Mean is %f.\n" % stats.mean(temp2018))
# standard deviation of the data
print("The Standard Deviation is %f.\n" %
      stats.stdev(temp2018))
# variance of the data
print("The Variance is %f.\n" % stats.variance(temp2018))
# skewness of the data
print("The Skewness is %f.\n" % skewness(temp2018))
# kurtosis of the data
print("The Kurtosis is %f.\n" % kurtosis(temp2018))
# median of the data
print("The Median is %f.\n" % stats.median(temp2018))
# percentiles
print("The 5th, 25th, 75th and 95th percentiles are:")
probs = [5, 25, 75, 95]
print([round(np.percentile(temp2018, p), 5) for p in probs])
print()

```

6316

6317 The meaning of these indices may be explained as follows. The mean is the simple
 6318 average of samples. The variance of climate data reflects the strength of variations
 6319 of a climate system and has units of the square of the data entries, such as [$^{\circ}$ C]².
 6320 You may have noticed the denominator $n - 1$ instead of n , which is for an estimate
 6321 of unbiased sample variance. The standard deviation describes the spread of the
 6322 sample entries and has the same units as the data. A large standard deviation
 6323 means that the samples have a broad spread.

6324 Skewness is a dimensionless quantity. It measures the asymmetry of sample data.
 6325 Zero skewness means a symmetric distribution about the sample mean. For example,
 6326 the skewness of a normal distribution is zero. Negative skewness denotes a skew to
 6327 the left, meaning the existence of a long tail on the left side of the distribution.
 6328 Positive skewness implies a long tail on the right side.

6329 The words “Kurtosis” and “kurtic” are Greek in origin and indicate peakedness.
 6330 Kurtosis is also dimensionless and indicates the degree of peakedness of a probability
 6331 distribution. The kurtosis of a normal distribution¹ is zero when 3 is subtracted as
 6332 in Eq. (B.7). Positive kurtosis means a high peak at the mean, thus the distribution
 6333 shape is slim and tall. This is referred to as leptokurtic. “Lepto” is Greek in origin
 6334 and means thin or fine. Negative kurtosis indicates a low peak at the mean, thus
 6335 the distribution shape is fat and short, referred to as platykurtic. “Platy” is also
 6336 Greek in origin and means flat or broad.

6337 For the 139 years of the NOAAGlobalTemp global average annual mean tem-
 6338 perature anomalies, mean is -0.1959°C , which means that the 1880-2018 average is
 6339 lower than the average during the climatology period 1971-2000. During the clima-

¹ Chapter 2 will have a detailed description of the normal distribution and other probabilistic distributions.

6340 tology period, the temperature anomaly average is approximately zero, and can be
 6341 computed by the R command `mean(temp2018[92:121])`.

6342 The variance of the data in the 139 years from 1880 to 2018 is $0.1055 [{}^{\circ}\text{C}]^2$,
 6343 and the corresponding standard derivation is $0.3248 [{}^{\circ}\text{C}]$. The skewness is 0.7743,
 6344 meaning skew to the right with a long tail on the right, thus with more extreme high
 6345 temperatures than extreme low temperatures, as shown by Fig. B.4. The kurtosis
 6346 is -0.2619, meaning the distribution is flatter than a normal distribution, as shown
 6347 in the histogram Fig. B.4.

6348 The median is -0.2744°C and is a number characterizing a set of samples such
 6349 that 50% of the sample values are less than the median, and another 50% are
 6350 greater than the median. To find the median, sort the samples from the smallest to
 6351 the largest. The median is then the sample value in the middle. If the number of
 6352 the samples is even, then the median is equal to the mean of the two middle sample
 6353 values.

6354 Quantiles are defined in a way similar to median by sorting. For example, 25-
 6355 percentile (also called the 25th percentile) -0.4120°C is a value such that 25% of
 6356 sample data are less than this value. By definition, 60-percentile is thus larger than
 6357 50-percentile. Here, percentile is a description of quantile relative to 100. Obviously,
 6358 100-percentile is the largest datum, and 0-percentile is the smallest one. Often, a
 6359 box plot is used to show the typical quantiles (see Fig. B.5).

6360 The 50-percentile (or 50th percentile) -0.2744°C is equal to the median. If the
 6361 distribution is symmetric, then the median is equal to mean. Otherwise these two
 6362 quantities are not equal. If the skew is to the right, then the mean is on the right of
 6363 the median: the mean is greater than the median. If the skew is to the left, then the
 6364 mean is on the left of the median: the mean is less than the median. Our 139 years
 6365 of temperature data are right skewed and have mean equal to -0.1959°C , greater
 6366 than their median equal to -0.2969°C .

6367 B.2 Commonly used climate statistical plots

6368 We will use the 139 years of NOAAGlobalTemp temperature data and R to illustrate
 6369 some commonly used statistical figures: histogram, boxplot, q-q plot, and linear
 6370 regression trend line.

6372 B.2.1 Histogram of a set of data

6373 A histogram of the NOAAGlobalTemp global average annual mean temperature
 6374 anomalies data from 1880-2018 is shown in Fig. B.4, which can be generated by the
 6375 following R code.

```
6376
6377 #R plot Fig. 1.4: Histogram and its fit
6378 par(mar=c(3.5,3.5,2.5,1), mgp=c(2,0.8,0))
6379 h <- hist(NOAAtemp[1:139, 2],
```

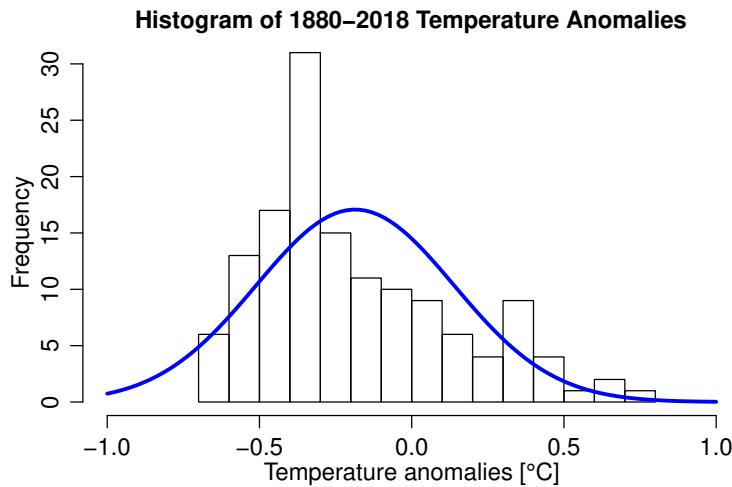


Fig. B.4 Histogram and its normal distribution fit of the global average annual mean temperature anomalies from 1880–2018. Each small interval in the horizontal coordinate is called a bin. The frequency in the vertical coordinate is the number of temperature anomalies in a given bin. For example, the frequency for the bin [-0.5, -0.4]°C is 17.

```

6380 main="Histogram of 1880–2018 Temperature Anomalies",
6381 xlab=expression(paste(
6382   "Temperature anomalies [", degree, "C]")),
6383 xlim=c(-1,1), ylim=c(0,30),
6384 breaks=10, cex.lab=1.2, cex.axis=1.2)
6385 xfit <- seq(-1, 1, length=100)
6386 areat <- sum((h$counts)*diff(h$breaks[1:2]))#Normalization area
6387 yfit <- areat*dnorm(xfit,
6388   mean=mean(NOAAtemp[1:139,2]),
6389   sd=sd(NOAAtemp[1:139,2]))
6390 #Plot the normal fit on the histogram
6391 lines(xfit, yfit, col="blue", lwd=3)

```

```

# Python plot Fig. 1.4: Histogram and its fit
mu, std = scistats.norm.fit(y)
# create an evenly spaced sequence of 100 points in [-1,1]
points = np.linspace(-1, 1, 100)
plt.hist(y, bins = 20, range=(-1,1), color ='white',
          edgecolor = 'k', density = True);
plt.plot(points, scistats.norm.pdf(points, mu, std),
          color = 'b', linewidth = 3)
plt.title(r"Histogram\u2014of\u20141880\u2014\u20142018\u2014Temperature\u2014Anomalies",
          pad = 20)
plt.xlabel("Temperature\u2014Anomalies\u2014[$\degree$C]", size = 25, labelpad = 20)
plt.ylabel("Frequency", size = 25, labelpad = 20)
plt.show()

```

6392

6393 The shape of the histogram agrees with the characteristics predicted by the
 6394 statistical indices in the previous sub-section:

- 6395 (i) The distribution is asymmetric and skewed to the right with skewness equal to
 6396 0.7743.
- 6397 (ii) The distribution is platykurtic with a kurtosis equal to -0.2619, i.e., it is flatter
 6398 than the standard normal distribution indicated by the blue curve.

6399
6400

B.2.2 Box plot

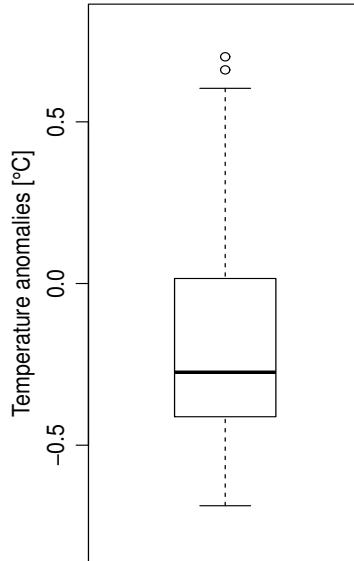
6401 Figure B.5 is the box plot of the 1880-2018 NOAAGlobalTemp global average annual
 6402 mean temperature anomaly data, and can be made from the following R command

```

#R plot Fig. 1.5: Box plot
boxplot(NOAAtemp[1:139, 2], ylim = c(-0.8, 0.8),
        ylab=expression(paste(
          "Temperature\u2014anomalies\u2014[", degree, "C]")),
        width=NULL, cex.lab=1.2, cex.axis=1.2)

```

6408 The rectangular box's mid line indicates the median, which is -0.2744°C . The
 6409 rectangular box's lower boundary is the first quartile, i.e., 25th percentile -0.4120°C .
 6410 The box's upper boundary is the third quartile, i.e., the 75th percentile 0.0155°C .
 6411 The box's height is the third quartile minus the first quartile, and is called the
 6412 interquartile range (IQR). The upper "whisker" is the third quartile plus 1.5 IQR.
 6413 The lower whisker is supposed to be at the first quartile minus 1.5 IQR. However,
 6414 this whisker would then be lower than the lower extreme. In this case, the lower
 6415 whisker takes the value of the lower extreme, which is -0.6872°C . The points outside
 6416 of the two whiskers are considered outliers. Our dataset has two outliers, which are
 6417 0.6607 and 0.7011°C , and are denoted by two small circles in the box plot. The
 6418 two outliers occurred in 2015 and 2016, respectively.



1

Fig. B.5 Box plot of the global average annual mean temperature anomalies from 1880-2018.

```
#Python plot Fig. 1.5: Box plot
medianprops = dict(linestyle='-', linewidth=2.5, color='k')
whiskerprops = dict(linestyle='--', linewidth=2.0, color='k')
plt.boxplot(y, medianprops = medianprops,
            whiskerprops = whiskerprops);
plt.title("Boxplot of 1880-2018 Temperature Anomalies",
          pad = 20)
plt.ylabel("Temperature Anomalies [$\degree$C]",
           size = 25, labelpad = 20)
y_ticks = np.linspace(-0.5,0.5,3)
plt.yticks(y_ticks)
plt.show()
```

6419

6420
6421

B.2.3 Q-Q plot

6422 Figure B.6 shows Quantile-Quantile (Q-Q) plots, also denoted by q-q plots,
6423 qq-plots, or QQ-plots.

6424 The function of a Q-Q plot is to compare the distribution of a given set of
6425 data with a specific reference distribution, such as a standard normal distribution
6426 with zero mean and standard deviation equal to one, denoted by $N(0, 1)$. A Q-Q
6427 plot lines up the percentiles of data on the vertical axis and the same number of

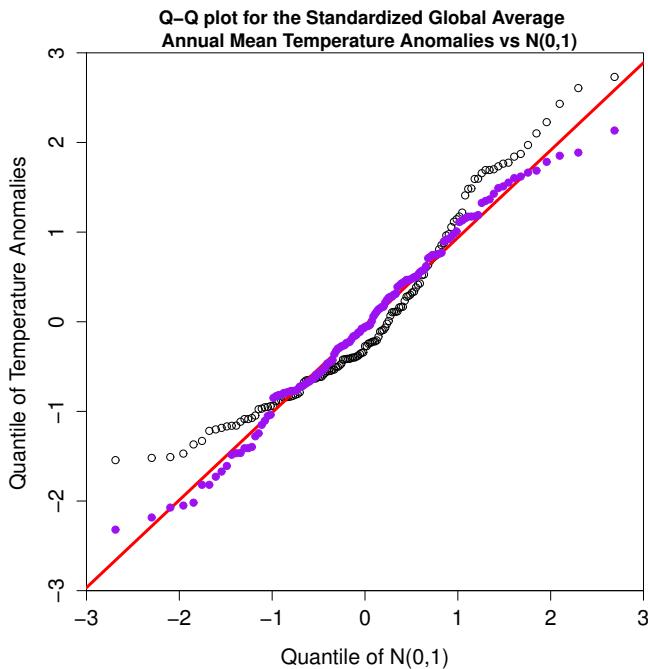


Fig. B.6 Black empty-circle points are the Q-Q plot of the standardized global average annual mean temperature anomalies vs. standard normal distribution. The purple points are the Q-Q plot for the data simulated by `rnorm(139)`. The red is the distribution reference line of $N(0, 1)$.

percentiles of the specific reference distribution on the horizontal axis. The pairs of the quantiles $(x_i, y_i), i = 1, 2, \dots, n$ determine the points on the Q-Q plot. Here, x_i and y_i correspond to the same cumulative percentage or probability p_i for both x and y variables, where p_i monotonically increases from approximately zero to one as i goes from 1 to n . A red Q-Q reference line is plotted as if the vertical axis values are also the quantiles of the given specific distribution. Thus, the Q-Q reference line should be diagonal.

The black empty circles in Fig.B.4 compare the quantiles of the standardized global average annual mean temperature anomalies marked on the vertical axis with those of the standard normal distribution marked on the horizontal axis. The standardized anomalies are equal to anomalies divided by the sample standard deviation. The purple dots shows a Q-Q plot of a set of 139 random numbers simulated by the standard normal distribution. As expected, the simulated points are located close to the red diagonal line, which is the distribution reference line of $N(0, 1)$. On the other hand, the temperature Q-Q plot shows a considerable degree of scattering of the points away from the reference line. We may intuitively conclude that the global average annual temperature anomalies from 1880 to 2018 are not exactly distributed according to a normal distribution, also known as Gaussian) distribu-

6446 tion. However, we may also conclude that the distribution of these temperatures
 6447 is not very far away from the normal distribution either, because the points on
 6448 the Q-Q plot are not very far away from the distribution reference line, and also
 6449 because even the simulated $N(0, 1)$ points are noticeably off the reference line for
 6450 the extremes.

6451 Figure B.6 can be generated by the following R code.

```

6452 # R plot Fig. 1.6: Q-Q plot for the standardized
6453 # global average annual mean temperature anomalies
6454 temp2018 <- NOAAtemp[1:139,2]
6455 tstand <- (temp2018 - mean(temp2018))/sd(temp2018)
6456 set.seed(101)
6457 qn <- rnorm(139) #simulate 139 points by N(0,1)
6458 qns <- sort(qn) # sort the points
6459 qq2 <- qqnorm(qns,col="blue", lwd = 2)
6460
6461 setEPS() #Automatically saves the eps file
6462 postscript("fig0106.eps", height=7, width=7)
6463 par(mar = c(4.5,5,2.5,1), xaxs = "i", yaxs = "i")
6464 qt = qnorm(tstand,
6465   main = "Q-Qplot for the Standardized Global Average
6466   Annual Mean Temperature Anomalies vs N(0,1)",
6467   ylab="Quantile of Temperature Anomalies",
6468   xlab="Quantile of N(0,1)",
6469   xlim=c(-3,3), ylim = c(-3,3),
6470   cex.lab = 1.3, cex.axis = 1.3)
6471 qqline(tstand, col = "red", lwd=3)
6472 points(qq2$x, qq2$y, pch = 19,
6473   col = "purple")
6474 dev.off()

```

6475 In the R code, we first standardize, also called normalize, the global average
 6476 annual mean temperature data by subtracting the data mean and dividing by the
 6477 data's standard deviation. Then, we use these 139 years of standardized global
 6478 average annual mean temperature anomalies to generate a Q-Q plot, which is shown
 6479 in Fig. B.6.

```

#Python plot Fig. 1.6: Q-Q plot for the standardized
# global average annual mean temperature anomalies
NOAAtemp = read_table(
    "data/aravg.ann.land_ocean.90S.90N.v4.0.1.201907.txt",
    header = None, delimiter = "\s+")
x = np.array(NOAAtemp.loc[0:138, 0])
y = np.array(NOAAtemp.loc[0:138, 1])
line = np.linspace(-3, 3, y.size)
tstand = np.sort((y - np.mean(y))/np.std(y))
# simulate 139 points following N(0,1)
qn = np.random.normal(size = y.size)
qns = np.sort(qn) # sort the points
qq2 = sm.qqplot(qns)
fig = plt.figure(figsize=(12,12)) # set up figure
sm.qqplot(tstand, color = "k", linewidth = 1)
# plot diagonal line
plt.plot(line, line, 'r-', linewidth = 3)
# Q-Q plot of standard normal simulations
plt.plot(line, qns, 'mo')
plt.tick_params(length=6, width=2, labelsize=20)
plt.title("Q-Q plot for the Standardized Global\nTemperature Anomalies vs N(0,1)", pad = 20)
plt.xlabel("Quantile of N(0,1)", size = 25, labelpad = 20)
plt.ylabel("Quantile of Temperature Anomalies", size = 25)
plt.show()

```

6480

6481
6482

B.2.4 Plot a linear trend line

6483 Climate data analysis often involves plotting a linear trend line for time series data,
 6484 such as the linear trend for the global average annual mean surface temperature
 6485 anomalies, shown in Fig. B.7. The R code for plotting a linear trend line of data
 6486 sequence y and time sequence t is `abline(y ~ t)`.

6487 Figure B.7 can be generated by the following computer code.

```

# R plot Fig. 1.7: Data line graph with a linear trend line
par(mar=c(3.5,3.5,2.5,1), mgp=c(2,0.8,0))
plot(NOAAtemp[1:139,1], NOAAtemp[1:139,2],
     type="l", col="brown", lwd=3,
     main="Global Land-Ocean Average Annual Mean
Surface Temperature Anomaly: 1880-2018",
     cex.lab=1.2,cex.axis=1.2,
     xlab="Year",
     ylab=expression(paste(
       "Temperature Anomaly [", degree,"C]")))
abline(lm(NOAAtemp[1:139,2] ~ NOAAtemp[1:139,1]),
       lwd=3, col="blue")
lm(NOAAtemp[1:139,2] ~ NOAAtemp[1:139,1])
# (Intercept) NOAAtemp[1:139, 1]
#-13.872921          0.007023
#Trend 0.7023 degC/100a
text(1930, 0.5,

```

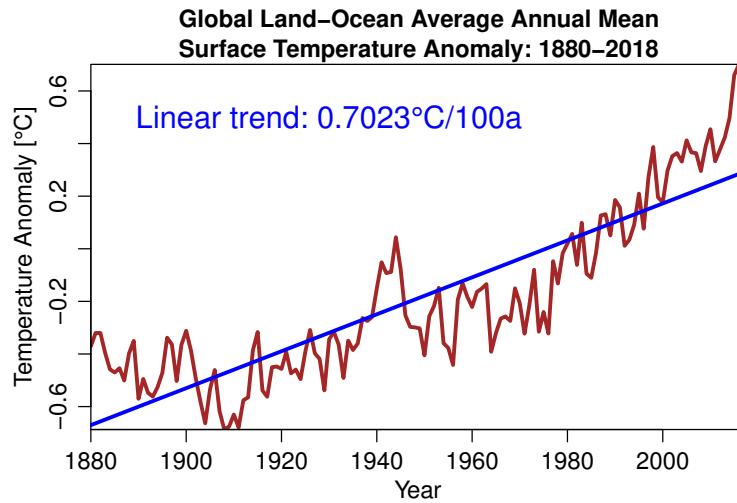


Fig. B.7 Linear trend line with the 1880–2018 global average annual mean surface temperature based on the NOAA GlobalTemp V4.0 dataset.

```

6506 expression(paste("Linear trend: 0.7023",
6507           degree, "C/100a")),
6508 cex = 1.5, col="blue")

# Python plot Fig. 1.7: Data line graph with a trend line
trend = np.array(np.polyfit(x, y, 1))
abline = trend[1] + x*trend[0]
plt.plot(x, y, 'k-',
          color = 'tab:brown', linewidth = 3);
plt.plot(x, abline, 'k-', color = 'b', linewidth = 3);
plt.title("Global Land–Ocean Average Annual Mean\nSurface Temperature Anomaly: 1880–2018", pad = 20);
plt.text(1880, 0.5, r"Linear trend: 0.7023$degree$C/100a",
         color= 'b', size = 28)
plt.xlabel("Year", size = 25, labelpad = 20)
plt.ylabel("Temperature Anomaly [$degree$C]", size = 25, labelpad = 20)

```

6510 B.3 Read netCDF data file and plot spatial data maps

6511

6512 B.3.1 Read netCDF data

6513

6514 Climate data are at spatiotemporal points, such as at the grid points on the Earth's
 6515 surface and at a sequence of time. NetCDF (Network Common Data Form) is a

6516 popular file format for modern climate data with spatial locations and temporal
 6517 records. The gridded NOAAGlobalTemp data has a netCDF version, and can be
 6518 downloaded from
 6519 <https://www.esrl.noaa.gov/psd/data/gridded/data.noaaglobaltemp.html>

6520 The data are written into a 3D array, with 2D latitude-longitude for space, and 1D
 6521 for time. R and Python can read and plot the netCDF data. We use the NOAA-
 6522 GlobalTemp as an example to illustrate the netCDF data reading and plotting.
 6523 Figure B.8 displays a temperature anomaly map for the entire globe for December
 6524 2015.

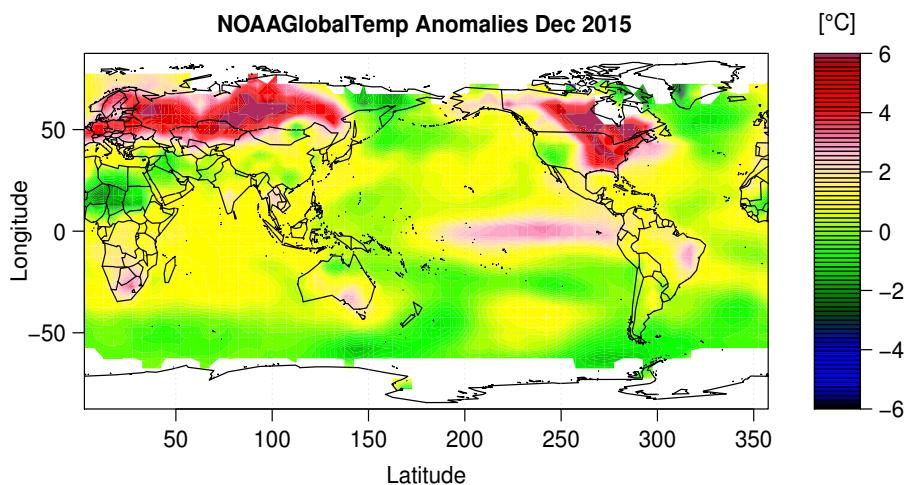


Fig. B.8 The surface temperature anomalies of December 2015 with respect to the 1971-2000 climatology. Data source: The NOAAGlobaTemp V4.0 gridded monthly data.

6525 Figure B.8 can be generated by the following computer code.

```
6526 # R read the netCDF data: NOAAGlobalTemp
6527 setwd("/Users/sshen/climstats")
6528 #install.packages("ncdf4")
6529 library(ncdf4)
6530 nc = ncdf4::nc_open("data/air.mon.anom.nc")
6531 nc # describes details of the dataset
6532 Lat <- ncvar_get(nc, "lat")
6533 Lat # latitude data
6534 #[1] -87.5 -82.5 -77.5 -72.5 -67.5 -62.5
6535 Lon <- ncvar_get(nc, "lon")
6536 Lon # longitude data
6537 #[1] 2.5 7.5 12.5 17.5 22.5 27.5
6538 Time <- ncvar_get(nc, "time")
6539 head(Time) # time data in Julian days
6540 #[1] 29219 29250 29279 29310 29340 29371
6541 library(chron) # convert Julian date to calendar date
6542 nc$dim$time$units # .nc base time for conversion
6543 #[1] "days since 1800-1-1 00:00:0.0"
6544 month.day.year(29219,c(month = 1, day = 1, year = 1800))
```

```

6545 #1880-01-01 # the beginning time of the dataset
6546 tail(Time)
6547 #[1] 79988 80019 80047 80078 80108 80139
6548 month.day.year(80139,c(month = 1, day = 1, year = 1800))
6549 #2019-06-01 # the end time of the dataset
6550
6551 # extract anomaly data in (lon, lat, time) coordinates
6552 NOAAgridT <- ncvar_get(nc, "air")
6553 dim(NOAAgridT) # dimensions of the data array
6554 #[1] 72 36 1674 #5-by-5, 1674 months from Jan 1880-Jun 2019

```

```

#Python read the netCDF data: NOAAGlobalTemp
import netCDF4 as nc # import netCDF data reading package
# go to the working directory
os.chdir('/Users/sshen/climstats')
# read a .nc file from the folder named "data"
nc = nc.Dataset('data/air.mon.anom.nc')
# get the detailed description of the dataset
print(nc)
# extract latitude, longitude, time and temperature
lon = nc.variables['lon'][:]
lat = nc.variables['lat'][:]
time = nc.variables['time'][:]
air = nc.variables['air']
# covert Julian date to calendar date
from netCDF4 import num2date
from datetime import datetime, date, timedelta
from matplotlib.dates import date2num, num2date
units = nc.variables['time'].units
print(units)
dtimes = num2date(time)

```

6555

B.3.2 Plot a spatial map of temperature

The NOAAGlobalTemp anomaly data on a 5-degree latitude-longitude grid for a given time can be represented by a color map. Figure B.8 shows the temperature anomaly map for December 2015, an El Niño month. The eastern tropical Pacific has positive anomalies, which is a typical El Niño signal. That particular month also exhibited a very large anomaly across Europe, and the eastern United States and Canada. The white areas over the high latitude regions lack data.

The figure can be generated by the following R code.

```

6565 #R plot Fig. 1.8: Dec 2015 surface temp anomalies map
6566 library(maps)# requires maps package
6567 mapmat=NOAAgridT[,,1632]
6568 # Julian date time 1632 corresponds to Dec 2015
6569 mapmat=pmax(pmin(mapmat,6),-6) # put values in [-6, 6]
6570 int=seq(-6, 6, length.out=81)
6571 rgbs.palette=colorRampPalette(c('black','blue',
6572 'darkgreen','green','yellow','pink','red','maroon'),
6573 interpolate='spline')

```

```

6574 par(mar=c(3.5, 4, 2.5, 1), mgp=c(2.3, 0.8, 0))
6575 filled.contour(Lon, Lat, mapmat,
6576   color.palette=rgb.palette, levels=int,
6577   plot.title=title(main="NOAAGlobalTempUAnomaliesUDecU2015",
6578     xlab="Latitude", ylab="Longitude", cex.lab=1.2),
6579   plot.axes={axis(1, cex.axis=1.2, las=1);
6580     axis(2, cex.axis=1.2, las=2); map('world2', add=TRUE); grid()},
6581   key.title=title(main=expression(paste("[", degree, "C]"))),
6582   key.axes={axis(4, cex.axis=1.2)})

```

```

# Python plot Fig. 1.8: Dec 2015 surface temp anomalies map
dpi = 100
fig = plt.figure(figsize = (1100/dpi, 1100/dpi), dpi = dpi)
ax = fig.add_axes([0.1, 0.1, 0.8, 0.9])
# create map
dmap = Basemap(projection = 'cyl', llcrnrlat = min(lat),
                urcrnrlat = max(lat), resolution = 'c',
                llcrnrlon = min(lon), urcrnrlon = max(lon))
# draw coastlines, state and country boundaries, edge of map
dmap.drawcoastlines()
dmap.drawstates()
dmap.drawcountries()
# convert latitude/longitude values to plot x/y values
x, y = dmap(*np.meshgrid(lon, lat))
# draw filled contours
cnplot = dmap.contourf(x, y, mapmat, clev, cmap=myColMap)
# tick marks
ax.set_xticks([0, 50, 100, 150, 200, 250, 300, 350])
ax.set_yticks([-50, 0, 50])
ax.tick_params(length=6, width=2, labelsize=20)
# add colorbar
# pad: distance between map and colorbar
cbar = dmap.colorbar(cnplot, pad = "4%", drawedges=True,
                     shrink=0.55, ticks = [-6,-4,-2,0,2,4,6])
# add colorbar title
cbar.ax.set_title('[\u00b0C]', size= 17, pad = 10)
cbar.ax.tick_params(labelsize = 15)
# add plot title
plt.title('NOAAGlobalTempUAnomaliesUDecU2015',
           size = 25, fontweight = "bold", pad = 15)
# label x and y
plt.xlabel('Longitude', size = 25, labelpad = 20)
plt.ylabel('Latitude', size = 25, labelpad = 10)
# display on screen
plt.show()

```

6583

B.3.3 Panoply plot of a spatial map of temperature

6584 You can also use the Panoply software package to plot the map (See Fig. B.9).
 6585 This is a very powerful data visualization tool developed by NASA specifically for

6588 displaying netCDF files. The software package is free and can be downloaded from
 6589 www.giss.nasa.gov/tools/panoply .

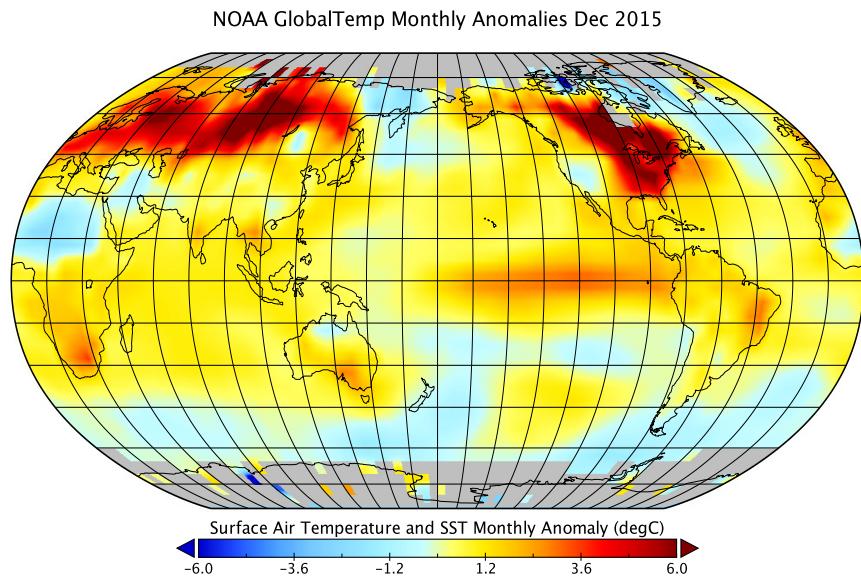


Fig. B.9 A Panoply plot of Robinson projection map for the surface temperature anomalies of December 2015.

6590 To make a Panoply plot, open Panoply and choose **Open** from the **File** menu.
 6591 Open dropdown which will allow you to go to the right directory to find the netCDF
 6592 file you wish to plot. In our case, the file is `air.mon.anom.nc`. Choose the climate
 6593 parameter `air`. Click on **Create Plot**. A map will show up. Then you have many
 6594 choices to modify the map, ranging from the data **Array**, **Scale**, and **Map**, ...,
 6595 and finally produce the figure. You can then tune the figure by choosing differ-
 6596 ent graphics parameters underneath the figure, such as **Array(s)** to choose which
 6597 month to plot, **Map** to choose the map projection types and the map layout options,
 6598 and **Labels** to type proper captions and labels.

6599 To learn more about Panoply, please use an online Panoply tutorial, such as the
 6600 NASA Panoply help page <https://www.giss.nasa.gov/tools/panoply/help/>.

6601 Compared with the R or Python map Fig. B.5, the visualization effect of the
 6602 Panoply map seems more appealing. However, R and Python have an advantage of
 6603 flexibility and can deal with all kinds of data. For example, the Plotly graphing li-
 6604 brary in R <https://plotly.com/r/> and in Python <https://plotly.com/python/>
 6605 can even make interactive and 3D graphics. You may find some high quality fig-
 6606 ures from the Intergovernmental Panel on Climate Change (IPCC) report (2021)
 6607 www.ipcc.ch and reproduce them using the computing tools described here.

B.4 1D-space-1D-time data and Hovmöller diagram

6609

6610 A very useful climate data visualization technique is the Hovmöller diagram. It
 6611 displays how a climate variable varies with respect to time along a given line section
 6612 of latitude, longitude or altitude. It usually has the abscissa for time and ordinate
 6613 for the line section. A Hovmöller diagram can conveniently show time evolution of
 6614 a spatial pattern, such as wave motion from south to north or from west to east.

6615 Figure B.10 is a Hovmöller diagram for the sea surface temperature (SST) anomalies
 6616 at a longitude equal to 240° , i.e., 120°W , in a latitude interval $[30^{\circ}\text{S}, 30^{\circ}\text{N}]$,
 6617 with time range from January 1989 to December 2018. When the red strips become
 6618 strong from the south to north, a strong El Niño occurs, such as those in the
 6619 1997-1998 and 2015-2016 winters.

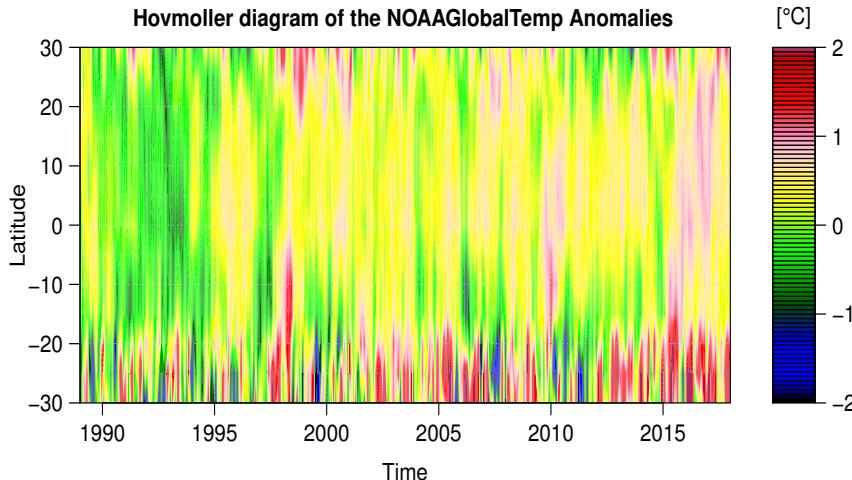


Fig. B.10 Hovmöller diagram for the gridded NOAA Global Temp monthly anomalies at longitude 120°W and a latitude interval $[30^{\circ}\text{S}, 30^{\circ}\text{N}]$.

6620 The Hovmöller diagram Fig. B.10 may be plotted by the following R code.

```
#R plot Fig. 1.10: Hovmoller diagram
library(maps)
mapmat=NOAAgridT[30,12:24,1309:1668]
#Longitude= 240 deg, Lat =[-30 30] deg
#Time=Jan 1989-Dec 2018: 30 years
mapmat=pmax(pmin(mapmat,2),-2) # put values in [-2,2]
par(mar=c(4,5,3,0))
int=seq(-2,2,length.out=81)
rgb.palette=colorRampPalette(c('black','blue',
'darkgreen','green','yellow','pink','red','maroon'),
interpolate='spline')
par(mar=c(3.5,3.5,2.5,1), mgp=c(2.4, 0.8, 0))
x = seq(1989, 2018, len=360)
y = seq(-30, 30, by=5)
```

```

6635 filled.contour(x, y, t(mapmat),
6636     color.palette=rgb.palette, levels=int,
6637     plot.title=title(main=
6638 "Hovmoller diagram of the NOAA Global Temp Anomalies",
6639         xlab="Time", ylab="Latitude", cex.lab=1.2),
6640     plot.axes={axis(1, cex.axis=1.2);
6641         axis(2, cex.axis=1.2);
6642         map('world2', add=TRUE); grid()},
6643     key.title=title(main =
6644         expression(paste("[", degree, "C]"))),
6645     key.axes={axis(4, cex.axis=1.2)})

```

```

# Python plot Fig. 1.10: Hovmoller diagram
mapmat2 = NOAAgridT[1308:1668,11:23,29]

# define values between -2 and 2
mapmat2 = np.array([[j if j < 2 else 2 for j in i]
                     for i in mapmat2])
mapmat2 = np.array([[j if j > -2 else -2 for j in i]
                     for i in mapmat2])
# find dimensions
mapmat2 = np.transpose(mapmat2)
print(mapmat2.shape)
lat3 = np.linspace(-30,30, 12)
print(lat3.shape)
time = np.linspace(1989, 2018, 360)
print(time.shape)

# plot functions
myColMap = LinearSegmentedColormap.from_list(
    name='my_list',
    colors=['black','blue','darkgreen','green','lime',
            'yellow', 'pink', 'red', 'maroon'], N=100)
clev2 = np.linspace(mapmat2.min(), mapmat2.max(), 501)
contf = plt.contourf(time, lat3, mapmat2,
                     clev2, cmap=myColMap);

plt.text(2019.2, 31.5,
         "[\u00b0C]", color='black', size = 23)
plt.title("Hovmoller diagram of the\nNOAA Global Temp Anomalies",
          fontweight = "bold", size = 25, pad = 20)
plt.xlabel("Time", size = 25, labelpad = 20)
plt.ylabel("Latitude", size = 25, labelpad = 12)
colbar = plt.colorbar(contf, drawedges=False,
                      ticks = [-2,-1,0,1,2])

```

6647 B.5 4D netCDF file and its map plotting

6648

6649 The 4D climate data means 3D spatial dimensions and 1D time. For example, the
 6650 NCEP Global Ocean Data Assimilation System (GODAS) monthly water tempera-
 6651 ture data are at 40 depth levels ranging from 5 meters to 4478 meters and at 1/3
 6652 degree latitude by 1 degree longitude horizontal resolution and are from January
 6653 1980. The NOAA-CIRES 20th Century Reanalysis (20CR) monthly air tempera-
 6654 ture data are at 24 different pressure levels ranging from 1000 mb to 10 mb and at
 6655 2-degree latitude and longitude horizontal resolution and are from January 1851.

6656 GODAS data can be downloaded from NOAA ESRL

6657 www.esrl.noaa.gov/psd/data/gridded/data.godas.html

6658 The data for each year is a netCDF file and is about 140MB. The following R code
 6659 can read the GODAS 2015 data into R.

```
6660 # R read a 4D netCDF file: lon, lat, level, time
6661 setwd("/Users/sshen/climstats")
6662 library(ncdf4)
6663 # read GODAS data 1-by-1 deg, 40 levels, Jan-Dec 2015
6664 nc=ncdf4::nc_open("data/godas2015.nc")
6665 nc
6666 Lat <- ncvar_get(nc, "lat")
6667 Lon <- ncvar_get(nc, "lon")
6668 Level <- ncvar_get(nc, "level")
6669 Time <- ncvar_get(nc, "time")
6670 head(Time)
6671 #[1] 78527 78558 78586 78617 78647 78678
6672 library(chron)
6673 month.day.year(78527,c(month = 1, day = 1, year = 1800))
6674 # 2015-01-01
6675 # potential temperature pottmp[lon, lat, level, time]
6676 godasT <- ncvar_get(nc, "pottmp")
6677 dim(godasT)
6678 #[1] 360 418 40 12,
6679 #i.e., 360 lon, 418 lat, 40 levels, 12 months=2015
6680 t(godasT[246:250, 209:210, 2, 12])
6681 #Dec level 2 (15-meter depth) water temperature [K] of
6682 #a few grid boxes over the eastern tropical Pacific
6683 # [,1] [,2] [,3] [,4] [,5]
6684 #[1,] 300.0655 299.9831 299.8793 299.7771 299.6641
6685 #[2,] 300.1845 300.1006 299.9998 299.9007 299.8045
```

```

# Python read a netCDF data file
import netCDF4 as nc
# go to your working directory
os.chdir('/Users/sshen/climstats')
# read a .nc file from the folder named "data"
nc1 = nc.Dataset('data/godas2015.nc')
# get the detailed description of the dataset
print(nc1)
# dimensions of the 2015 pottem data array
godasT = nc1.variables['pottmp'][:]
print(godasT.shape)
# (12, 40, 418, 360)

```

6686

Figure B.11 shows the 2015 annual mean water temperature at 195 meters depth based on the GODAS data. At this depth level, the equatorial upwelling appears: The deep ocean cooler water in the equatorial region upwells and makes the equatorial water cool. The equatorial water is not the hottest anymore at this level, and is cooler than the water in some subtropical regions as shown in Fig. B.11.

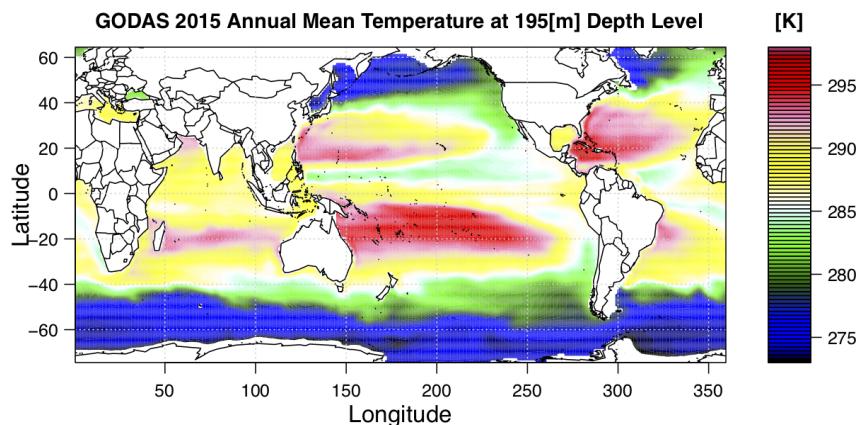


Fig. B.11 The 2015 annual mean water temperature at 195 meters depth based on GODAS data.

Figure B.11 can be generated by the following R code.

```

# R plot Fig. 1.11: The ocean potential temperature
# the 20th layer from surface: 195 meters depth
# compute 2015 annual mean temeprature at 20th layer
library(maps)
climmat=matrix(0,nrow=360,ncol=418)
sdmat=matrix(0,nrow=360,ncol=418)
Jmon<-1:12
for (i in 1:360){
  for (j in 1:418){
    climmat[i,j] = mean(godasT[i,j,20,Jmon]);
    sdmat[i,j]=sd(godasT[i,j,20,Jmon])
}

```

```

6704         }
6705     }
6706     int=seq(273,298,length.out=81)
6707     rgb.palette=colorRampPalette(c('black','blue',
6708       'darkgreen','green','white','yellow',
6709       'pink','red','maroon'), interpolate='spline')
6710     par(mar=c(3.5, 3.5, 2.5, 0), mgp=c(2, 0.8, 0))
6711     filled.contour(Lon, Lat, climmat,
6712       color.palette=rgb.palette, levels=int,
6713       plot.title=title(main=
6714         "GODAS\u20142015\u2014Annual\u2014Mean\u2014Temperature\u2014at\u2014195[m]\u2014Depth\u2014Level",
6715         xlab="Longitude", ylab="Latitude",
6716         cex.lab=1.3, cex.axis=1.3),
6717       plot.axes={axis(1); axis(2); map('world2', add=TRUE); grid()},
6718       key.title=title(main=" [K]"))

```

```

# Python plot Fig. 1.11: A spatial map from 4D data
climmat = np.zeros((360, 418))
for i in range(360):
    for j in range(418):
        climmat[i,j] = np.mean(godasT[:, 20, j, i])
climmat = np.transpose(climmat)
lat4 = np.linspace(-75, 65, 418)
long = np.linspace(0, 360, 360)

myColMap = LinearSegmentedColormap.from_list(
    name='my_list',
    colors=['black','blue','darkgreen','green',
            'white','yellow','pink','red','maroon'],
    N=100)
plt.figure(figsize=(18,12));
clev3 = np.arange(godasT.min(), godasT.max(), 0.1)
contf = plt.contourf(long, lat4, climmat,
                     clev3, cmap=myColMap);
plt.text(382, 66, "[K]", fontsize=23, color='black')
plt.tick_params(length=6, width=2, labelsize=20)
m = Basemap(projection='cyl', llcrnrlon=0,
             urcrnrlon=360, resolution='l', fix_aspect=False,
             suppress_ticks=False, llcrnrlat=-75, urcrnrlat=65)
m.drawcoastlines(linewidth=1)
plt.title("GODAS\u20142015\u2014Annual\u2014Mean\u2014Temperature\u2014\n\u2014at\u2014195[m]\u2014Depth\u2014Level",
           size = 25, fontweight = "bold", pad = 20)
plt.xlabel("Longitude", size = 25, labelpad = 20)
plt.ylabel("Latitude", size = 25, labelpad = 15)
plt.tick_params(length=6, width=2, labelsize=30)
colbar = plt.colorbar(contf, drawedges=False,
                      ticks = [275,280,285,290,295])

```

6720

B.6 Paraview and 4DVD

6721

6722 Besides using R, Python and Panoply to plot climate data, other software packages
 6723 may also be used to visualize data for some specific purposes, such as Paraview for
 6724 3D visualization and 4DVD for fast climate diagnostics and data delivery.

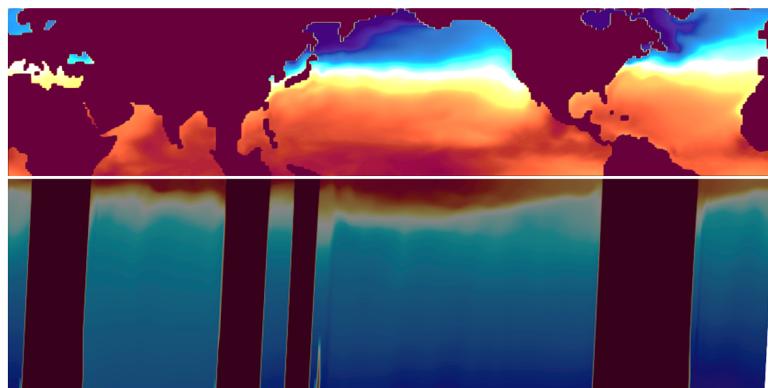
6725

6726

B.6.1 Paraview

6727 Paraview is an open-source data visualization software package, available at
 6728 www.paraview.org/download. There are many online tutorials, such as
 6729 www.paraview.org/tutorials.

6730 You may use the software ParaView to plot the GODAS data in a 3D view as
 6731 shown in Fig. B.12 for the December 2015 water temperature.

**Fig. B.12**

A 3D view of the December 2015 annual mean water temperature based on the GODAS data: The surface of the Northern Hemisphere and the cross-sectional map along the equator from surface to the 2,000 meters depth. The dark color indicates land. The red, yellow and blue colors indicate temperature.

6732

6733

B.6.2 4DVD

6734 4DVD (4-dimensional visual delivery of big climate data) is a fast data visualization
 6735 and delivery software system www.4dvd.org. It optimally harnesses the power of
 6736 distributed computing, database and data storage to allow a large number of general
 6737 public users to quickly visualize climate data. For example, teachers and students
 6738 can use 4DVD to visualize climate model data in classrooms and download the
 6739 visualized data instantly. The 4DVD website has a tutorial for users.

6740 Here, we provide an example of 4DVD visualization of the NOAA-CIRES 20th
 6741 Century Reanalysis climate model data for atmosphere. 4DVD can display a tem-
 6742 perature map at a given time and given pressure level as shown in Fig. B.13 for

6743 January 1851 and 750 millibar pressure level, or one can obtain a time series of
 6744 the monthly air temperature for a specific grid point from January 1851. In fact,
 6745 4DVD can show multiple time series for the same latitude-longitude location but
 6746 at different pressure levels. The 4DVD not only allows a user to visualize the data,
 6747 but also download the data for the figure shown in the 4DVD system. In this sense,
 6748 4DVD is like a machine that plays data, while the regular DVD player machine,
 6749 popular for about 30 years since 1980s, play DVD discs for music and movies.

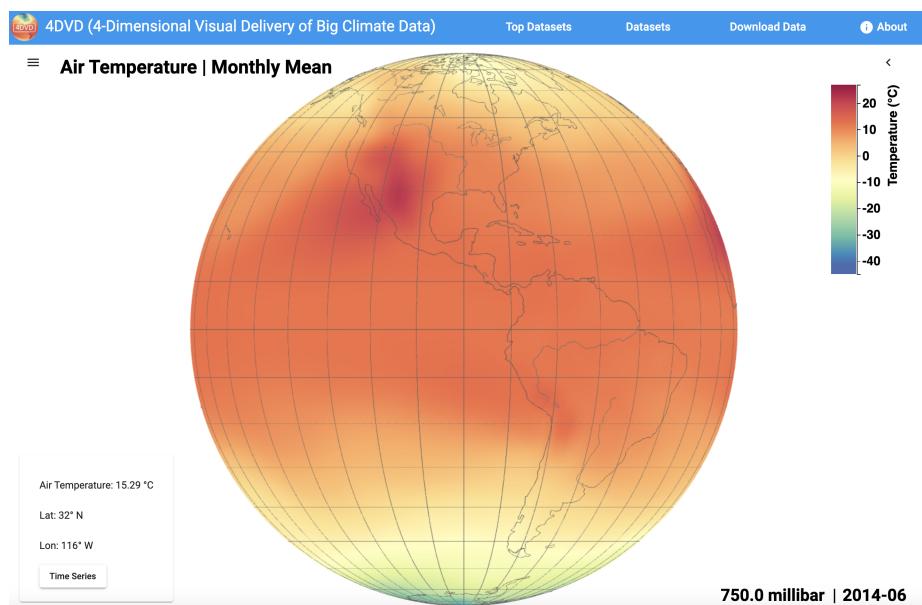


Fig. B.13 The 4DVD screenshot for the the NOAA-CIRES 20th Century Reanalysis temperature at 750 millibar height level for June 2014.

B.6.3 Other online climate data visualization tools

6750
 6751 Besides R, Python, Panoply, ParaView, and 4DVD, there are many other data
 6752 visualization and delivery software systems. A few popular and free ones are listed
 6753 as follows.

6754
 6755 Nullschool <https://nullschool.net> is a beautiful data visualizer of wind, ocean
 6756 flows, and many climate parameters. It is supported by the data from a global
 6757 numerical weather prediction system, named the Global Forecast System (GFS)
 6758 run by the United States' National Weather Service (NWS).

6759 Ventusky www.ventusky.com has both website and smartphone app. It has an
 6760 attractive and user-friendly interface that allows users to get digital weather informa-
 6761 tion instantly around the globe.

6762 Climate Reanalyzer climatedataanalyser.org is a comprehensive tool for climate
 6763 data plotting and download. It has a user-friendly interface for reanalysis and his-

6764 torical station data. The data can be exported in either CSV (comma-separated
 6765 values) format or JSON (JavaScript object notation) format.

6766 Google Earth Engine earthengine.google.com provides visualization tools to-
 6767 gether with a huge multi-petabytes storage of climate data. Its modeling and satel-
 6768 lite data sources are from multiple countries.

6769 Giovanni giovanni.gsfc.nasa.gov is an online climate data plotting tool with
 6770 an interface. It allows to download the plotted figures in different formats, such as
 6771 png. It is supported by various kinds of NASA climate datasets.

6772 Climate Engine climateengine.org is a web application for plotting climate
 6773 and remote sensing data. Similar to Giovanni, it also has a tabular interface for a
 6774 user to customized her data and maps.

6775 NOAA Climate at a Glance www.ncdc.noaa.gov/cag is a data visualization tool
 6776 mainly for visualizing the observed climate data over the United States. It has
 6777 functions for both spatial maps and historical time series.

6778 Web-based Reanalyses Intercomparison Tools (WRIT) psl.noaa.gov/data/writ.
 6779 Similar to Giovanni, WRIT also has an interface table for a user to enter plot pa-
 6780 rameters. The plot (in postscript format) and its data (in netCDF format) can be
 6781 downloaded. WRIT is designed for the data from climate reanalysis models.

6782 **B.6.4 Use ChatGPT as a study assistant**

6783 ChatGPT (Generative Pre-trained Transformer) is an artificial intelligence (AI)
 6784 system that can help us develop computer codes, write essays, draft course syllabi,
 6785 and more. For example, in your box of the ChatGPT interface, you can enter “Write
 6786 a Python code to plot the cosine function as a blue thick dashed curve.” ChatGPT
 6787 will give you the following Python code:

```
6784 import numpy as np
6785 import matplotlib.pyplot as plt
6786
6787 x = np.linspace(-np.pi, np.pi, 100)
6788 y = np.cos(x)
6789
6790 plt.plot(x, y, linestyle='--', linewidth=3, color='blue')
6791 plt.xlabel('x')
6792 plt.ylabel('cos(x)')
6793 plt.title('Plot of the cosine function')
6794 plt.show()
```

6790 You can copy the above code and paste it to your Jupyter Notebook cell, then run
 6791 the code to produce the curve of the cosine function in the interval $[-\pi, \pi]$. With
 6792 this sample and according to your needs, you can modify the code for different
 6793 functions, intervals, colors, thickness, title, axis labels, and more.

6794 You can also ask ChatGPT to work on a dataset, such as “Write an R code
 6795 to plot the NOAAGlobalTempts.csv data and its trend line.” ChatGPT will give

6796 you an R code. You copy and paste the code to your RStudio. Set your R code in
 6797 the correct directory that has the NOAAGlobalTemp annual time series data file,
 6798 `NOAAGlobalTemps.csv`. Run your R code to plot the data and the linear trend
 6799 line. Of course, you can modify the code according to your requirements. The R
 6800 code is as follows.

```
6801 setwd('/Users/sshen/climmath/data')
6802 library(tidyverse)
6803
6804 # Load the data
6805 data <- read.csv("NOAAGlobalTemps.csv")
6806
6807 # Plot the data
6808 ggplot(data, aes(x = Year, y = Tempts)) +
6809   geom_line() +
6810   geom_smooth(method = "lm", se = FALSE) +
6811   ggtitle("NOAA Global Temperature Data") +
6812   xlab("Year") +
6813   ylab("Temperature (Celsius)")
```

6814 Another example is that you enter “Write an essay on global warming.” ChatGPT
 6815 will write an essay for you. Most likely, it does not reflect your idea, but it will
 6816 give you a few tips, such as the definition of global warming, evidences and data,
 6817 consequences, main causes, how to slow down the warming, and conclusion. Again,
 6818 you can use this ChatGPT essay as a sample and modify the essay using your own
 6819 data and ideas, such as green technologies, solar farms, microclimate resources for
 6820 smart cities, trend of the oil business, lifestyles, AI in self-driving electric cars, and
 6821 education.

6822 ChatGPT was released by the OpenAI Lab in November 2022. You may use this
 6823 tool as your assistant of learning and working. It is only your assistant. It is not
 6824 you! You can efficiently use it to give you hints and inspire your ideas. You still
 6825 have to produce your own work.

6826 B.7 Chapter summary

6828 This chapter has provided a brief introduction to useful statistical concepts and
 6829 methods for climate science and has included the following materials.

6830 (i) Formulas to compute the most commonly used statistical indices:

- 6831 • Mean as a simple average, median as a datum whose value is in the
 6832 middle of the sorted data sequence, i.e., the median is larger than 50%
 6833 of the data and smaller than the remaining 50%,
- 6834 • Standard deviation as a measure of the width of the probability distri-
 6835 bution,
- 6836 • Variance as the square of the standard deviation,
- 6837 • Skewness as a measure of the degree of asymmetry in distribution, and

- 6838 • Kurtosis as a measure of the peakedness of the data distribution com-
6839 pared to that of the normal distribution.

6840 (ii) The commonly used statistical plots:

- 6841 • Histogram for displaying the probability distribution of data,
6842 • Linear regression line for providing a linear model for data,
6843 • Box plot for quantifying the probability distribution of data, and
6844 • Q-Q plot for checking whether the data are normally distributed.

6845 (iii) Read and plot the netCDF file for plotting a 2D map by R and Python.

6846 Other data visualization tools, such as Panoply, Paraview, 4DVD, Plotly,
6847 and Nullschool, were briefly introduced. The online tools like 4DVD and
6848 Nullschool may be used for classroom teaching and learning.

6849 Computer codes and climate data examples are given to demonstrate these con-
6850 cepts and the use of the relevant tools and formulas. With this background plus
6851 some R or Python programming skill, you will have sufficient knowledge to meet
6852 the needs of the basic statistical analysis for climate data.

References and Further Readings

- 6854 [1] Hennemuth, B., Bender, S., K. Bulow, N. Dreier, E. Keup-Thiel, O. Kruger,
 6855 C. Mudersbach, C. Radermacher, and R. Schoetter, 2013: *Statistical methods*
 6856 *for the analysis of simulated and observed climate data, applied in projects and*
 6857 *institutions dealing with climate change impact and adaptation.* CSC Report 13,
 6858 Climate Service Center, Germany, 135pp. URL:
 6859 [http://www.climate-service-center.de/imperia/md/content/csc/
 6860 projekte/csc-report13_englisch_final-mit_umschlag.pdf](http://www.climate-service-center.de/imperia/md/content/csc/projekte/csc-report13_englisch_final-mit_umschlag.pdf)

This free statistics recipe book outlines numerous methods for climate data analysis, which are collected and edited by climate science professionals and have examples of real climate data.

- 6861
 6862 [2] IPCC, 2021: AR6 Climate Change 2021: The Physical Science Basis. Contribution
 6863 of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change [Masson-Delmotte, V., P. Zhai, A. Pirani, S.
 6864 L. Connors, C. Pan, S. Berger, N. Caud, Y. Chen, L. Goldfarb, M. I. Gomis, M.
 6865 Huang, K. Leitzell, E. Lonnoy, J. B. R. Matthews, T. K. Maycock, T. Waterfield,
 6866 O. Yeleki, R. Yu and B. Zhou (eds.)]. Cambridge University Press.

This is the famous IPCC report for free at the IPCC website, <https://www.ipcc.ch/report/ar6/wg1/>. It includes many high quality figures plotted from climate data

- 6868
 6869 [3] NCEI, 2021: *Anomalies vs. Temperature* Last accessed on 12 May 2021.

6870 [https://www.ncdc.noaa.gov/monitoring-references/dyk/
 6871 anomalies-vs-temperature](https://www.ncdc.noaa.gov/monitoring-references/dyk/anomalies-vs-temperature)

6872

This is an excellent site for education and gives an easy-to-understand justification of the use of anomalies. From this site, you can find many other educational resources for climate science, such as the concise description of climate extreme index and global precipitation percentile maps.

6873

- 6874 [4] Shen, S.S.P, and R.C.J. Somerville, 2019: Climate Mathematics: Theory and
 6875 Applications. Cambridge University Press, 391pp.

This book attempts to modernize the mathematical education for undergraduate students majoring in atmospheric and oceanic sciences, or related fields. The book integrates six traditional mathematics courses (i.e., Calculus I, II, and III, Linear Algebra, Statistics, and Computer Programming) into a single course named Climate Mathematics. The course uses real climate data and can be taught in three semesters for freshmen and sophomores, or in one semester as an upper level or graduate class. Computer codes of R and Python and other learning resources are available at the book website www.climatemathematics.org.

- 6876 [5] Smith, T.M., R.W. Reynolds, T.C. Peterson, and J. Lawrimore, 2008: Improvements to NOAA's historical merged landcean surface temperatures analysis
 6877 (1880006). Journal of Climate, 21, 2283296, doi:10.1175/2007JCLI2100.1.
 6878

These authors are among the main contributors who reconstructed the sea surface temperature (SST). They began their endeavor in the early 1990s.

- 6880 [6] Vose, R.S., D. Arndt, V.F. Banzon, D.R. Easterling, B. Gleason, B. Huang,
 6881 E. Kearns, J.H. Lawrimore, M.J. Menne, T.C. Peterson, R.W. Reynolds, T.M.
 6882 Smith, C.N. Williams, Jr., and D.L. Wuertz, 2012: NOAA's merged land-ocean
 6883 surface temperature analysis. Bulletin of the American Meteorological Society,
 6884 93, 1677685.
 6885

These authors are experts on the reconstruction of both SST and the land surface air temperature, and have published many papers in data quality control and uncertainty quantifications.

6887 Exercises

- 6888 **B.1** The NOAA Merged Land Ocean Global Surface Temperature Analysis (NOAA-
 6889 GlobalTemp) V5 includes the global land-spatial average annual mean tem-
 6890 perature anomalies as shown below

| | | | | | |
|------------|-----------|----------|----------|----------|----------|
| 6892 1880 | -0.843351 | 0.031336 | 0.009789 | 0.000850 | 0.020698 |
| 6893 1881 | -0.778600 | 0.031363 | 0.009789 | 0.000877 | 0.020698 |
| 6894 1882 | -0.802413 | 0.031384 | 0.009789 | 0.000897 | 0.020698 |
| 6895 | | | | | |

- 6896 The first column is for time in years, and the second is the temperature anomalies in [Kelvin]. Columns 3-6 are data errors. This data file for the land temperature can be downloaded from the NOAAGlobalTemp website
 6897 www.ncdc.noaa.gov/data/noaa-global-surface-temperature/v5/access/timeseries
 6898
 6899
 6900 The data file named `aravg.ann.land.90S.90N.v5.0.0.202104.asc.txt` is
 6901 also included in `data.zip` that can be downloaded from the book website
 6902 www.climatestatistics.org
 6903 (a) Plot the anomalies against time from 1880 to 2020 using the point-line
 6904 graph like Fig. B.1.
 6905 (b) Plot the anomalies against time from 1880 to 2020 using the staircase
 6906 chart like Fig. B.2.
 6907 (c) Plot the anomalies against time from 1880 to 2020 using the color bar
 6908 chart like Fig. B.3.
- 6909 **B.2** NOAAGlobalTemp V5 also has the global **ocean**-spatial average annual mean
 6910 temperature anomalies contained in a data file named
 6911 `aravg.ann.locean.90S.90N.v5.0.0.202104.asc.txt`
 6912 For this dataset, please plot the anomaly data in the same three styles (a),
 6913 (b) and (c) as in the previous problem. You may download the data file from
 6914 the NOAAGlobalTemp V5 website or from `data.zip` for this book.
- 6915 **B.3** NOAAGlobalTemp V5 also has the global land-spatial average monthly mean
 6916 temperature anomalies from January 1880 to April 2021. The data file is
 6917 named
 6918 `aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt`
 6919 For this dataset, please plot the January anomaly data from January 1880 to
 6920 January 2021 in the same three styles (a) - (c) as in the previous problem.
- 6921 **B.4** For the monthly global land data of NOAAGlobalTemp V5 in the above prob-
 6922 lem,
 6923 (a) For the January data from 1880 to 2020, compute the following statistical
 6924 indices: Mean, standard deviation, variance, skewness, kurtosis, max, 75%-
 6925 percentile, median, 25%-percentile, and min.
 6926 (b) Do the same for February, March, ..., December.
 6927 (c) Aggregate all the results in (a) and (b) into a 10×12 matrix with the 10
 6928 statistical indices in rows and the 12 months in columns. Add proper column
 6929 names, such as Jan Feb Mar ..., and also row names.
 6930 (d) Use 100-300 words to describe the differences of the statistical indices for
 6931 different months.
- 6932 **B.5** For the monthly data CRUTEM4 (Climate Research Unit surface air tempera-
 6933 ture anomalies) in 4DVD, download the historical time series data for January
 6934 of a location of your interest, and compute the following statistical indices:
 6935 Mean, standard deviation, variance, skewness, kurtosis, max, 75%-percentile,
 6936 median, 25%-percentile, and min.
- 6937 **B.6** Do the same as the previous problem but for July. Comment on the differences
 6938 between the statistical indices in January and July.
- 6939 **B.7** (a) Plot a histogram of the January global land temperature anomalies using

- 6940 the data file in the above problem, i.e.,
6941 `aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt`
6942 (b) Do the same for July.
6943 (c) Use 100-200 words to describe the comparison of the two histograms.
- B.8** (a) Plot a box plot for the January global land temperature anomalies
6944 `aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt`
6945 (b) Do the same but for the July data.
6946 (c) Put the two box plots next to each other in the same figure.
6947 (d) Use 100-200 words to describe the comparison of the two box plots.
- B.9** (a) Use the monthly data in the above problem to generate twelve (12) Q-Q plots relative to the standard normal distribution for every month from January to December.
6949 (b) From the 12 Q-Q plots, discuss which month's temperature anomalies are
6950 the closest to the normal distribution.
6951
- B.10** (a) Using the monthly data in the above problem
6952 `aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt`
6953 compute the linear trend of January temperature anomalies from 1880 to
6954 2020. Output your result in the unit: [°C/century].
6955 (b) Repeat (a) for each of the other eleven months.
6956 (c) Generate a list for the twelve linear trends.
6957 (d) Use 100-200 words to discuss the numerical values of the list.
- B.11** Plot the December 1997 surface temperature anomalies using the NOAA-
6958 GlobalTemp V5 data on a $5^\circ \times 5^\circ$ grid in the netCDF format:
- 6963 `NOAAGlobalTemp_v5.0.0_gridded_s188001_e202104_c20210509T133251.nc`
- 6964 The data can be downloaded from NOAAGlobalTemp V5 or extracted from
6965 `data.zip` for this book at the website www.climatestatistics.org
- B.12** Use Panoply to make the same December 1997 surface temperature anomalies
6966 plot as the above problem, but with the *Robinson* map projection.
- B.13** Use Panoply to make the same December 1997 surface temperature anomalies
6967 plot as the above problem, but with the *Mercator* map projection.
- B.14** Use Panoply to make the same December 1997 surface temperature anomalies
6968 plot as the above problem, but with the *Orthographic* map projection.
- B.15** Use Panoply to make the same December 1997 surface temperature anomalies
6969 plot as the above problem, but with the *Stereographic (Two-Hemisphere)* map
6970 projection.
- B.16** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp
6971 monthly anomalies at longitude 150°W and a latitude interval $[50^\circ\text{S}, 50^\circ\text{N}]$
6972 from January 1989 to December 2018 (i.e., 240 months).
- B.17** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp
6973 monthly anomalies at longitude 140°W and a latitude interval $[40^\circ\text{S}, 40^\circ\text{N}]$
6974 from January 1971 to December 2000 (i.e., 360 months).

- 6981 **B.18** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp
6982 monthly anomalies on the equator with longitude from 160°E to 90°W] from
6983 January 1971 to December 2000 (i.e., 360 months).
- 6984 **B.19** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp
6985 monthly anomalies at latitude 5°S with longitude from 160°E to 90°W] from
6986 January 1971 to December 2000 (i.e., 360 months).
- 6987 **B.20** Use R or Python to plot four December 2015 potential temperature maps
6988 at the depth layers 5, 25, 105 and 195 meters based on the GODAS data.
6989 You can use the netCDF data file `godas2015.nc` in `data.zip` for this book
6990 or download the netCDF GODAS data from a NOAA website, such as
6991 <https://www.esrl.noaa.gov/psd/data/gridded/data.godas.html>
- 6992 **B.21** Use Panoply to plot the same maps as the previous problem but with the
6993 Robinson projection.
- 6994 **B.22** Use Plotly in R or Python to plot four December 2015 potential temperature
6995 maps at the depth layers 5, 25, 105 and 195 meters based on the GODAS data.
6996 Try to place the four maps together to make a 3D visualization.
- 6997 **B.23** Plot four June 2015 potential temperature maps at the depth layers 5, 25,
6998 105 and 195 meters using the same netCDF GODAS data file as the above
6999 problem.
- 7000 **B.24** Use Plotly in R or Python to plot four June 2015 potential temperature maps
7001 at the depth layers 5, 25, 105 and 195 meters based on the GODAS data. Try
7002 to place the four maps together to make a 3D visualization.
- 7003 **B.25** Plot three cross-sectional maps of the December 2015 potential temperature
7004 at 10°S, equator, and 10°N using the same netCDF GODAS data file as the
7005 above problem. Each map is on a 40×360 grid, with 40 depth levels and 1-deg
7006 longitude resolution for the equator.
- 7007 **B.26** Do the same as the above but for the June 2015 GODAS potential temper-
7008 ature data.
- 7009 **B.27** Use R or Python a cross-sectional map of the December 2015 potential tem-
7010 perature along a meridional line at longitude at 170°E using the same netCDF
7011 GODAS data file as the above problem.
- 7012 **B.28** Use R or Python to plot a cross-sectional map of the June 2015 potential
7013 temperature along a meridional line at longitude at 170°E using the same
7014 netCDF GODAS data file as the above problem.

Index