

<sup>1</sup>  
<sup>2</sup>  
**Interactive Linear Algebra with  
R and Python**

<sup>3</sup>  
**SAMUEL S.P. SHEN**

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

<sup>7</sup> August 2023 Version for SDSU Course Math 524



# Contents

<i>Preface</i>	<i>page</i>	vii
<i>Acknowledgements</i>		xii
<b>1 Basic Matrix Operations</b>	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 Mass balance for chemical equations in marine chemistry	25	
1.8 Multivariate linear regression using matrix notations	26	
1.9 Chapter summary	28	
References and Further Readings	31	
Exercises	32	
<b>2 Matrix Theory and Visualization</b>	35	
2.1 Matrix Definitions	35	
2.2 Fundamental Properties of Matrices	38	
2.3 Some basic concepts and theories of linear algebra	43	
2.3.1 Linear equations	43	

44	2.3.2 Linear transformations	44
45	2.3.3 Linear independence	45
46	2.3.4 Determinants	46
47	2.3.5 Rank of a matrix	47
48	2.4 Eigenvectors and eigenvalues	50
49	2.4.1 Definition of eigenvectors and eigenvalues	50
50	2.4.2 Properties of eigenvectors and eigenvalues for a symmetric matrix	54
51		
52	2.5 Hadamard and Other Matrix Multiplications	56
53	2.5.1 Hadamard Product of Two Matrices of the Same Dimensions	57
54	2.5.2 Jordan Product of Two Matrices of the Same Dimensions	57
55	2.5.3 Commutator of Two Matrices of the Same Dimensions	57
56	2.5.4 Outer Product of Two Vectors and Two Matrices	58
57	2.5.5 Kronecker Product of Two Matrices of the Same Dimensions	60
58	2.6 Direct Sum of Two Matrices	61
59	2.7 Visualization of Eigenvalues and Eigenvectors for a Covariance Matrix	61
60		
61	2.8 Singular Value Decomposition	62
62	2.8.1 SVD formula and a simple SVD example	62
63	2.9 SVD for the standardized sea level pressure data of Tahiti and Darwin	68
64		
65	2.10 Chapter summary	70
66	<i>References and Further Readings</i>	72
67	Exercises	73
68	<b>3 Matrix Applications to Machine Learning</b>	77
69	3.1 K-means clustering	77
70	3.1.1 K-means setup and trivial examples	78
71	3.1.2 A K-means algorithm	84
72	3.1.3 K-means clustering for the daily Miami weather data	86
73	3.2 Support vector machine	96
74	3.2.1 SVM for a system of three points labeled in two categories	101
75	3.2.2 SVM mathematical formulation for a system of many points in two categories	105
76		
77	3.3 Random forest method for classification and regression	111
78	3.3.1 RF flower classification for a benchmark iris dataset	111
79	3.3.2 RF regression for the daily ozone data of New York City	118
80	3.3.3 What does a decision tree look like?	122
81	3.4 Neural network and deep learning	125
82	3.4.1 An NN model for an automated decision system	125
83	3.4.2 An NN prediction of iris species	132
84	3.5 Chapter summary	136
85	<i>References and Further Readings</i>	137
86	Exercises	138

87	<b>4 Matrix Applications to Regression Models</b>	141
88	4.1 Simple linear regression	141
89	4.1.1 Temperature lapse rate and an approximately linear model	141
90	4.1.2 Assumptions and formula derivations of the single variate	
91	linear regression	145
92	4.1.3 Statistics of slope and intercept: Distributions, confidence	
93	intervals, and inference	156
94	4.2 Multiple linear regression	169
95	4.2.1 Calculating the Colorado TLR when taking location	
96	coordinates into account	169
97	4.2.2 Formulas for estimating parameters in the multiple linear	
98	regression	172
99	4.3 Nonlinear fittings using the multiple linear regression	174
100	4.3.1 Diagnostics of linear regression: An example of global	
101	temperature	174
102	4.3.2 Fit a third order polynomial	179
103	4.4 Linear Regression by Weighted Least Squares	182
104	4.5 Chapter summary	183
105	<i>References and Further Readings</i>	184
106	Exercises	185
107	<b>Appendix A A Tutorial of R and RStudio</b>	188
108	A.1 Download and install R and R-Studio	188
109	A.2 R Tutorial	189
110	A.2.1 R as a smart calculator	190
111	A.2.2 Define a sequence in R	191
112	A.2.3 Define a function in R	192
113	A.2.4 Plot with R	192
114	A.2.5 Symbolic calculations by R	193
115	A.2.6 Vectors and matrices	194
116	A.2.7 Simple statistics by R	197
117	A.3 Online Tutorials	198
118	A.3.1 YouTube tutorial: for true beginners	199
119	A.3.2 YouTube tutorial: for some basic statistical summaries	199
120	A.3.3 YouTube tutorial: Input data by reading a csv file into R	199
121	A.4 Chapter summary	201
122	<i>References and Further Readings</i>	203
123	Exercises	204
124	<b>Appendix B Visualization of Matrices</b>	207
125	<i>References and Further Readings</i>	239
126	<i>Index</i>	245

<sup>127</sup> Give the pupils something to do, not something to learn; and doing is  
<sup>128</sup> of such a nature as to demand thinking; learning naturally results.

<sup>129</sup>

<sup>130</sup>

— John Dewey

## Preface

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

In 1965 the Committee on the Undergraduate Program in Mathematics, Mathematical Association of America (MAA), outlined the following topics for a stand-alone linear algebra course: Linear systems, matrices, vectors, linear transformations, unitary geometry with characteristic values.

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

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

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

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

**167    What are the distinguished features of this book?**

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

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

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

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

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

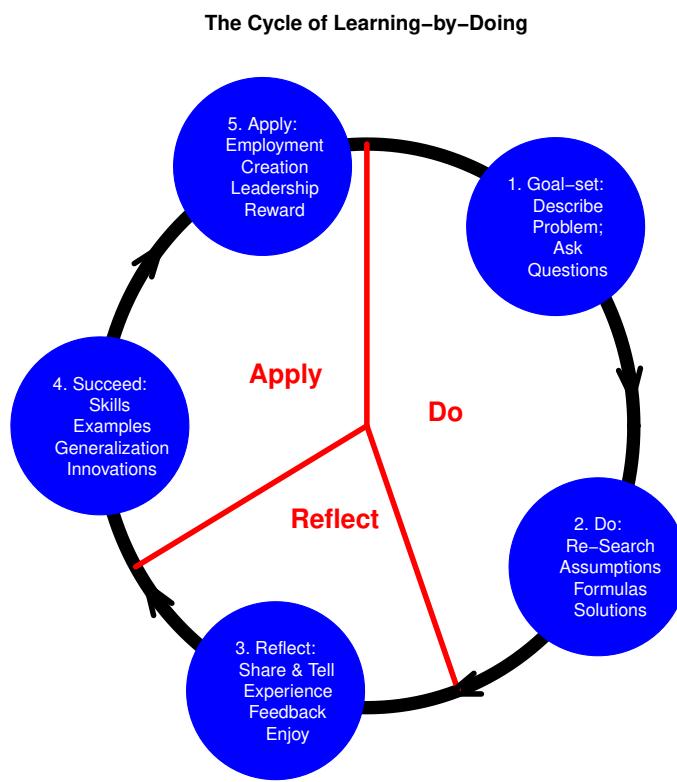
209    Another feature of this book is to minimize the use of mathematical symbols  
210    and notations. We understand that professional mathematicians do not mind fancy

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

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

We follow the education theory of *learning-by-doing*, which in our case means using your computer and our R or Python code to interact with our book and modifying our computer codes and websites to analyze your own data and generate corresponding figures. Learning-by-doing is the core methodology of the progressive education of John Dewey (1859-1952), an American philosopher and educator. Dewey felt that the experience of students and teachers together yields extra value for both. Instructors are not just to lecture and project authority, instead they are to collaborate with students and guide to students to gain experience of solving problems of their interest. Although Dewey's education theory was established initially for school children, we feel that the same is applicable to undergraduate and graduate students. Our way of learning-by-doing is to enable students to use R or Python code and other resources in the book and its website to reproduce the figures and numerical results in the book. Further, students can modify the computer code and solve their own problems, such as visualizing the health data or weather data in a similar format and of their own interest, or analyzing their own data using a similar or a modified method. Thus, the audience interaction with the book is the main innovative feature of our book, which allows the audience to gain experience of practicing, thinking, applying, and consequently understanding. The 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 Confucius were approximately 2,500 years apart, they shared a similar philosophy of learning-by-doing.

As illustrated in Fig. 0.1, our pedagogy has three stages: Do, reflect, and apply. The author has systematically practiced this pedagogy since 2015 when he created and taught the course *Climate Mathematics* at the Scripps Institution of Oceanography. Usually, he presents a question or a problem at the beginning of a class. Then he asks students to orally ask the same question, or describe the same problem, or answer his question in their own words. For example, why the tiny amount of carbon dioxide in the atmosphere is important for climate change? Next, he and students search and re-search for data and literature, work on the observed carbon dioxide data at Mauna Loa using computer data visualization to see the 30% increase of carbon dioxide in the atmosphere since 1958, and discuss the structure of greenhouse gasses whose molecules have three or more atoms. Next, he encourages 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 to 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

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

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

286

Samuel S. P. Shen, San Diego, California, August 2023

287 **References:**

- 288 1. Beaumont, R. A., 1965. Linear Algebra. Harcourt, Brace & World, Inc., New  
289 York.
- 290 2. Birkhoff, G. and Mac Lane, S., 2017. A survey of modern algebra. 3rd edition,  
291 Macmillan Company, New York, 1965.
- 292 3. Curtis, C.W., 1963. Linear algebra: an introductory approach. Allyn and Bacon,  
293 Inc., Boston.
- 294 4. Garcia, S.R. and Horn, R.A., 2023. Matrix Mathematics: A Second Course in  
295 Linear Algebra. 2nd edition, Cambridge University Press.
- 296 5. Grothendieck, A., Dieudonn, J. and Dieudonn, J., 1964. Elments de gomtrie  
297 algrique. Hermann, Paris.
- 298 6. Halmos, P.R., 1958. Finite-dimensional vector spaces. 2nd edition, D. Van Nos-  
299 trand Company Inc., Princeton.
- 300 7. Strang, G., 2016. Introduction to linear algebra. 5th edition, Wellesley-Cambridge  
301 Press.
- 302 8. Tucker, A., 1993. The growing importance of linear algebra in undergraduate  
303 mathematics. The college mathematics journal, 24(1), pp.3-9.

## Acknowledgements

<sup>305</sup> We thank the friendly editorial team of Cambridge University Press (CUP) for their  
<sup>306</sup> professional and high quality work.

309

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

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

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

327

## 1.1 Matrix as a data array

328

**329** In general, a matrix is a rectangular array of numbers or symbols or expressions,  
**330** which are called elements, arranged in rows or columns. A table such as that shown  
**331** in Fig. 1.1 is a matrix, consisting of  $N$  rows and  $Y$  columns of numbers. Figure  
**332** 1.1 shows the 10 ten-year annual precipitation anomalies from the year 1900 to  
**333** 1909 for the 15 five-degree latitude-longitude boxes centered at  $2.5^{\circ}E$  longitude for  
**334** different latitudes over the Northern Hemisphere ranging from latitudes  $2.5^{\circ}N$  to  
**335**  $72.5^{\circ}N$ . The matrix shown in Fig. 1.1 thus has 15 rows ( $N = 15$ ), and 10 columns  
**336** ( $Y = 10$ ). An anomaly of a climate parameter is defined as its actual value minus  
**337** its normal value that is an average of 30 or more years.

**338** Many other types of climate data can also be represented as matrices (which  
**339** is the plural of matrix). Precipitation data [units: mm/day] at multiple stations  
**340** and multiple days can also form a matrix, normally with stations [each marked  
**341** by a station identifier, or station ID] represented in rows, and time [units: day]  
**342** 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

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

### 1.2.1 Matrix equality, addition and subtraction

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

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

372 or

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

373 or simply

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

374 where  $a_{ij}, i = 1, \dots, m, j = 1, \dots, n$  are the  $mn$  elements of the rectangular matrix  
 375  $A$ , and  $m \times n$  is often called the size or order of a matrix. We say that  $A$  is an  $m$   
 376 times  $n$  matrix, or an  $m$ -by- $n$  matrix. The elements of the matrix shown in Fig. 1.1  
 377 are the precipitation anomaly data,  $m = 15$ , and  $n = 10$ . Therefore, Fig. 1.1 is a  
 378 15 times 10 matrix.

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

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

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

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

<sup>393</sup> relative humidity, etc, while the first columns represent time from 1:00 (i.e., 1:00  
<sup>394</sup> AM) to 24:00 (i.e., 12:00 AM).

<sup>395</sup> Matrix subtraction is defined in a similar way:

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

<sup>396</sup>

---

<sup>397</sup> **Example 1.1** Matrix subtraction:

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

<sup>398</sup>

---

<sup>400</sup> The R code for the above matrix subtraction can be written as follows:

```
401 matrix(c(1,1,1,-1), nrow=2) - matrix(c(1,3,2,0), nrow=2)
402 # [,1] [,2]
403 #[1,] 0 -1
404 #[2,] -2 -1
```

<sup>405</sup> One can use matrix subtraction to quantify the differences of climate between  
<sup>406</sup> two space-time domains when the climate data for each domain are in a matrix  
<sup>407</sup> form.

---

## <sup>408</sup> 1.2.2 Matrix multiplication

---

<sup>410</sup> While the matrix addition and subtraction are similar to those of scalars, the matrix  
<sup>411</sup> multiplication has several distinct properties.

### <sup>412</sup> 1.2.2.1 A row vector times a column vector

<sup>413</sup> The single-column n-row matrix is often called an n-dimensional vector, or an n-  
<sup>414</sup> dimensional column vector. Similarly, one can define an n-dimensional row vector  
<sup>415</sup> as a single-row n-column matrix.

<sup>416</sup> A row vector of  $n$  elements times a column vector of the same number of elements  
<sup>417</sup> is equal to a scalar, which is the sum of the products of each pair of corresponding  
<sup>418</sup> elements:

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

<sup>419</sup> This is also called the dot product of two vectors  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  and  $\mathbf{b} =$   
<sup>420</sup>  $(b_1, b_2, \dots, b_n)$ , denoted by

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n. \quad (1.8)$$

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

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

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

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

428

---

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

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

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

```
435 #R code for dot product
436 #install.packages('geometry')
437 library(geometry)
438 a = c(1, 1)
439 b = c(1, 0)
440 # Calculating dot product using dot()
441 dot(a, b)

442 #Another way to compute dot product is to
443 a%*%b
444 # [1,]   [,1]
445 # [1,]     1
446
447 #Calculate the angle between two vectors
448 am = norm(a, type="2")
449 am
450 am
451 # [1] 1.414214
452 bm = norm(b, type="2")
453 bm
454 # [1] 1
455 angleab = acos(dot(a, b)/(am*bm))*180/pi
456 angleab
457 # [1] 45 degrees
```

458

---

459

### 1.2.2.2 A scalar times a matrix

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

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

462 The product is again a matrix of the same size. A physically meaningful product  
 463 requires that the dimension of the scalar and the dimensions of the matrix elements  
 464 are compatible. For example, if the matrix is the price data of many products, and  
 465 if the scalar is the number of sales of each product, then the scalar times the matrix  
 466 gives the revenue matrix of all the products.

467

---

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

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

469 is

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

```
470 #R code for a scalar times a matrix
471 A = matrix(c(1,1,1,-1), ncol = 2, byrow=T)
472 A
473 #      [,1] [,2]
474 #[1,]    1    1
475 #[2,]    1   -1
476 3*A
477 #      [,1] [,2]
478 #[1,]    3    3
479 #[2,]    3   -3
```

480

---

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

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

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

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

490

491 **Example 1.4** Matrix multiplications:

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

492 and

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

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

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

496 In the same way, one can verify every element of the above multiplication results.  
 497 These matrix products can be computed by the following computer code.

```
498 #R code for matrix multiplication by command %*%
499 A = matrix(c(1,1,1,-1), nrow=2)
500 B = matrix(c(1,2,3,4), nrow=2)
501 A %*% B
502 #      [,1] [,2]
503 # [1,]    3    7
504 # [2,]   -1   -1
505 B %*% A
506 #      [,1] [,2]
507 # [1,]    4   -2
508 # [2,]    6   -2
```

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

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

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

514

---

#### 515 1.2.2.4 Transpose matrix, and diagonal, identity and zero 516 matrices

517 Although a space-time climate data matrix often has its rows to mark spatial loca-  
 518 tions and columns to mark time, the row and column roles may need to exchange  
 519 for some applications, which use rows to mark time and columns to mark spatial  
 520 locations. This operation of exchanging rows and columns is called matrix trans-  
 521 pose. The transposed matrix of  $A$  is denoted by  $A^t$ . The columns of  $A^t$  are the rows  
 522 of  $A$ .

523

**Example 1.5**

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

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

```
525 #R code for matrix transpose
526 A = matrix(c(1,2,3,4), ncol =2, byrow=T)
527 A
528 #      [,1] [,2]
529 #[1,]    1    2
530 #[2,]    3    4
531 t(A)
532 #      [,1] [,2]
533 #[1,]    1    3
534 #[2,]    2    4
```

535 It is obvious that

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

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

537

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

538

539

---

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

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

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

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

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

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

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

```

549 #Generate a diagonal matrix
550 D = diag(c(2, 1, -3))
551 round(D, digits = 0)
552 # [,1] [,2] [,3]
553 #[1,] 2 0 0
554 #[2,] 0 1 0
555 #[3,] 0 0 -3
556
557 #Generate an 3-dimensional identity matrix
558 I = diag(3)
559 I
560 # [,1] [,2] [,3]
561 #[1,] 1 0 0
562 #[2,] 0 1 0
563 #[3,] 0 0 1
564
565 #Generate a 2-by-3 zero matrix
566 M = matrix(0, 2, 3)
567 M
568 # [,1] [,2] [,3]
569 #[1,] 0 0 0
570 #[2,] 0 0 0

```

### 571 1.2.2.5 Matrix division and inverse

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

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

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

576 Matrix division is defined in the same way:

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

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

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

578 where  $I$  is the identity matrix.

---

579  
 580 **Example 1.6** The R command to find the inverse of matrix  $A$  is `solve(A)` as  
 581 shown by a numerical example below.

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

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

```

583 #R code to compute the inverse of a matrix
584 A = matrix(c(1,1,1,-1), nrow=2)
585 invA = solve(A)

```

```

586 invA
587 #      [,1] [,2]
588 #[1,]  0.5  0.5
589 #[2,]  0.5 -0.5
590
591 #Verify the inverse
592 invA %*% A
593 #      [,1] [,2]
594 #[1,]    1    0
595 #[2,]    0    1
596

```

---

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

## 602     **1.3 Some Useful Operations for Matrices in Data** 603     **Science but Not Covered in Traditional Linear Algebra** 604     **Books**

---

### 606     **1.3.1 Delete a Row or Column or Both; a Sub-Matrix**

---

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

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

```

617 #R code to delete rows or/and columns
618 A = matrix(1:9, nrow = 3)
619 A
620 #      [,1] [,2] [,3]
621 #[1,]    1    4    7
622 #[2,]    2    5    8
623 #[3,]    3    6    9
624
625 A[-3,]
626 #      [,1] [,2] [,3]
627 #[1,]    1    4    7
628 #[2,]    2    5    8

```

```

629 A[, -1]
630 # [,1] [,2]
631 #[1,] 4 7
632 #[2,] 5 8
633 #[3,] 6 9
635 A[-3, -1]
636 # [,1] [,2]
637 #[1,] 4 7
638 #[2,] 5 8
639
640 A[1:2, 2:3] #Sub-matrix
641 # [,1] [,2]
642 #[1,] 4 7
643 #[2,] 5 8
644
645 A[1:2, ] #keep all the columns
646 # [,1] [,2] [,3]
647 #[1,] 1 4 7
648 #[2,] 2 5 8

```

### 1.3.2 Insert Rows or/and Columns to a Matrix

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

```

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

```

```

682 cbind(Abc, A)#stack two matrices
683 #          bc
684 #[1,] 1 3 7 1 3
685 #[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.

```

691 #Row or column statistics
692 A = matrix(1:6, nrow = 2)
693 rowSums(A)
694 #[1] 9 12
695 rowMeans(A)
696 #[1] 3 4
697 rowCumsums(A) #cumulative sum
698 # [,1] [,2] [,3]
699 #[1,] 1 4 9
700 #[2,] 2 6 12
701 colMeans(A)
702 #[1] 1.5 3.5 5.5
703 library(matrixStats) #SD needs the library
704 rowSds(A)
705 #[1] 2 2
706 colSds(A)
707 #[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.

```

713 #Sweep a matrix by a vector using subtraction: R code
714 A = matrix(1:6, nrow=2, byrow = TRUE)
715 A
716 # [,1] [,2] [,3]
717 #[1,] 1 2 3
718 #[2,] 4 5 6
719 u = 1:3
720 Br = sweep(A, 2, u) #2 means sweep every row
721 Br
722 # [,1] [,2] [,3]
723 #[1,] 0 0 0
724 #[2,] 3 3 3
725 v= 1:2
726 Bc = sweep(A, 1, v) #1 means sweep every column
727 Bc
728 # [,1] [,2] [,3]
729 #[1,] 0 1 2

```

```

730 # [2,]    2    3    4
731
732 c = colMeans(A) #means of each column
733 sweep(A, 2, c) #compute the anomaly data matrix
734 #[1,] -1.5 -1.5 -1.5
735 #[2,]  1.5  1.5  1.5
736
737 sin(A)#function operation on each matrix element
738 #      [,1]     [,2]     [,3]
739 #[1,]  0.8414710  0.9092974  0.1411200
740 #[2,] -0.7568025 -0.9589243 -0.2794155
741
742 A^2 #not equal to A%*%A
743 #      [,1] [,2] [,3]
744 #[1,]    1    4    9
745 #[2,]   16   25   36

```

746     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  
747     may use R command `c = colMeans(A)`, and the anomaly data matrix would be  
748     the result of row sweeping by subtraction.

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

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

### 1.3.4.2 Multiplication sweeping

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

```

759 #R sweep a matrix by a vector using multiplication
760 A = matrix(1:6, nrow=2, byrow = TRUE)
761 w = 1:2
762 w*A
763 #      [,1] [,2] [,3]
764 #[1,]    1    2    3
765 #[2,]    8   10   12
766 A*w # yields the same result as w*A
767 w3 = 1:3
768 #sweep each row by transposing A
769 t(w3*t(A))
770 #      [,1] [,2] [,3]
771 #[1,]    1    4    9
772 #[2,]    4   10   18
773 w3*A ##multiplication sweep by row-dimensions not matching
774 #      [,1] [,2] [,3]
775 #[1,]    1    6    6
776 #[2,]    8    5   18
777
778 A/w #sweeping by division

```

```

779 # [,1] [,2] [,3]
780 #[1,] 1 2.0 3
781 #[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.

```

799 #Conversions between a Vector and a Matrix
800 v = c(60, 58, 67, 70, 55, 53)
801 M = matrix(v, nrow = 2) #from vector to matrix
802 M
803 # [,1] [,2] [,3]
804 #[1,] 60 67 55
805 #[2,] 58 70 53
806 c(M) #from matrix to vector by column
807 #[1] 60 58 67 70 55 53
808 c(t(M)) #from matrix to vector by row
809 #[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.

```

822 #R code to reduce the dimension of an nD array
823 x <- array(1:(2*3*4), dim=c(2,3,4))
824 dim(x)
825 #[1] 2 3 4
826 x #a stack of four 2-by-3 matrices
827 #, , 1 #the first of the 3rd dim index
828
829 #[,1] [,2] [,3]
830 #[1,]    1     3     5
831 #[2,]    2     4     6
832 # ...
833 #install.packages('R.utils')
834 library(R.utils)
835 #flat all the other dim except the 3rd one
836 #flat the 1st and 2nd dim
837 y <- wrap(x, map=list(3, NA))
838 dim(y)
839 #[1] 4 6
840 y
841 #[,1] [,2] [,3] [,4] [,5] [,6]
842 #[1,]    1     2     3     4     5     6
843 #[2,]    7     8     9    10    11    12
844 #[3,]   13    14    15    16    17    18
845 #[4,]   19    20    21    22    23    24
846
847 #back to the original 3D array
848 array(t(y), dim = c(2,3,4))

```

849

## 1.4 A Set of Linear Equations

850

851 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  
 852 engineering. Finding numerical solutions for a climate model based on partial differential equations will usually involve solving large systems of linear equations.  
 853

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

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

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

860 where  $x_1$  and  $x_2$  stand for the brothers' ages.

861 The matrix form of these equations would be

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

863 which involves three matrices:

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

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

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

871 Although one can easily guess that the solution to the above simple matrix equa-  
872 tion (1.28) is  $x_1 = 12$  and  $x_2 = 8$ , a more general method for computing the solution  
873 may be using the R code shown below:

```
874 #Solve linear equations
875 A = matrix(c(1,1,1,-1), nrow = 2)
876 b = c(20,4)
877 solve(A, b)
878 #[1] 12 8 #This is the result x1=12, and x2=8.
```

879 This type of R program can solve more complicated systems of linear equations,  
880 such as a system with 1,000 unknowns rather than two unknowns, as in this exam-  
881 ple.

882 The solution may be represented as

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

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

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

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

## 886 1.5 Eigenvalues and eigenvectors of a square space 887 matrix

---

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

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

---

895      **1.5.1 Matrices of data anomalies, standardized anomalies,  
896      covariance, and correlation**  
897

---

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

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

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

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

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

910      where

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

911      is the climate data mean of station  $i$  ( $i = 1, 2, \dots, N$ ). The mean is also called  
912      climatology or climate normal of the climate parameter. Thus, the temporal mean  
913      of the anomaly data for each station should be zero. However, in many climate  
914      applications, the temporal data are not continuous and were missing in the earlier  
915      period or in the middle. It is impossible to compute the temporal mean for the entire  
916      period of length  $Y$ . Instead, the climatology is computed using a sub-interval of data  
917      period, whose length is smaller than  $Y$ , for example, using 30 years from 1971 to  
918      2000. The 30-year climatology has almost become a recent community standard  
919      because of the period's best space-time coverage of observations. Consequently,  
920      the temporal mean of the anomaly data based on the 30-year climatology is often  
921      non-zero.

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

928      When the anomaly matrix is normalized by the standard deviation of each sta-  
929      tion, then the result is called the standardized anomaly matrix:

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

930 where

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

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

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

943 Then, the covariance and correlation matrices can be estimated by the following  
 944 formulas

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

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

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

---

947  
 948 **Example 1.7** This example's data matrix  $A$  has  $N = 2$  and  $Y = 3$ . The data's  
 949 covariance matrix is `covm`.

```
950 #Spatial covariance matrix
951 dat = matrix(c(0,-1,1,2,3,4), nrow=3)
952 dat
953 colMeans(dat)
954 A = sweep(dat, 2, colMeans(dat))
955 A
956 # [,1] [,2]
957 #[1,] 0 -1
958 #[2,] -1 0
959 #[3,] 1 1
960 covm=(1/(dim(A)[2]))*A%*%t(A)
961 covm #is the covariance matrix.
962 # [,1] [,2] [,3]
963 #[1,] 0.5 0.0 -0.5
964 #[2,] 0.0 0.5 -0.5
965 #[3,] -0.5 -0.5 1.0
```

---

### 1.5.2 Eigenvectors and their corresponding eigenvalues

---

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

```
970 u = c(1, 1, 0)
971 v = covm %*% u
972 v
973 #      [,1]
974 #[1,]  0.5
975 #[2,]  0.5
976 #[3,] -1.0
977 #u and v are in different directions
```

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

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

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

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

```
989 #Eigenvectors of a covariance matrix
990 ew = eigen(covm)
991 ew
992 #$values
993 #[1] 1.500000e+00 5.000000e-01 1.332268e-15
994
995 #$vectors
996 #      [,1]      [,2]      [,3]
997 #[1,] -0.4082483 -7.071068e-01 0.5773503
998 #[2,] -0.4082483  7.071068e-01 0.5773503
999 #[3,]  0.8164966  8.881784e-16 0.5773503
1000
1001 #Verify the eigenvectors and eigenvalues
1002 covm %*% ew$vectors[,1]/ew$values[1]
1003 #      [,1]
1004 #[1,] -0.4082483
1005 #[2,] -0.4082483
1006 #[3,]  0.8164966
1007 #This is the first eigenvector
1008
1009 w = ew$vectors[,1] # is an eigenvector
```

1010 A  $3 \times 3$  covariance matrix has three eigenvalues, and three eigenvectors  $(\lambda_1, w_1)$   
1011  $(\lambda_2, w_2)$ , and  $(\lambda_3, w_3)$ . These are called eigenpairs.

1012 In this example,

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

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

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

1013 The eigenvectors are frequently called modes, or empirical orthogonal functions  
1014 (EOFs) in climate science. The term EOF was coined by Edward Lorenz (1917-  
1015 2008) in his 1956 paper on statistical weather prediction (Lorenz 1956).

1016 Each eigenvalue is equal to the variance of the data projection on the corre-  
1017 sponding eigenvector, and is thus positive. The sum of all the eigenvalues of such  
1018 an  $N \times N$  matrix represents the total variance of the climate system observed at  
1019 these  $N$  stations. The first eigenvalue  $\lambda_1$  is the largest, corresponding to the largest  
1020 spatial variability of the climate field under study. The eigenvalues sizes follow the  
1021 order  $\lambda_1 \geq \lambda_2 \geq \dots$ .

1022 However, in carrying out an analysis of climate data, one can often find the  
1023 important patterns as eigenvectors more directly from the anomaly data matrix  
1024  $A$  without computing the covariance matrix  $C$  explicitly. This is known as the  
1025 singular value decomposition (SVD) approach (Golub and Reinsch 1970), which we  
1026 discuss next. It separates the space-time anomaly data into a space part, a time  
1027 part, and what we may think of as a variation in energy part. This mathematical  
1028 method of space-time decomposition is universally applicable to any data that we  
1029 sample in space and time, and it can often help to develop physical insight and  
1030 scientific understanding of the phenomena and their properties as contained in  
1031 the observational data. Efficient computing methods of SVD have been extensively  
1032 researched and developed since the 1960s. Gene H. Golub (1932-2007) was a leading  
1033 researcher in this effort and is remembered as “Professor SVD” by his Stanford  
1034 colleagues and the world mathematics community.

## 1035 1.6 An SVD representation model for space-time data

1036

### 1037 1.6.1 Space-Time Data Matrix and Its decomposition

1038 We encounter space-time data every day, a simple example being the air tempera-  
1039 ture at different locations at different times. If you take a plane to travel from San  
1040 Diego to New York, you may experience the temperature at San Diego in the morn-

1041

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  $\delta_{lk}$  is the Kronecker delta equal to

1071 zero when  $k \neq l$  and one when  $k = l$ . The columns of  $V$  are also a set of orthonormal  
 1072 vectors known as temporal eigenvectors.

1073 Usually, the elements of the  $U$  and  $V$  matrices are unitless (i.e., dimensionless),  
 1074 and the unit of the  $D$  elements is the same as the elements of the data matrix. For  
 1075 example, if  $A$  is a space-time precipitation data matrix with a unit [mm/day], then  
 1076 the dimension of the  $D$  elements is also [mm/day].

1077

### 1.6.2 SVD of the Space-Time Anomaly Data Matrix and 1078 Eigenvalue Problem of a Covariance Matrix

1079

---

1080

1081 **Example 1.8** SVD for a  $2 \times 3$  matrix.

```
1082 #Develop a 2-by-3 space-time data matrix for SVD
1083 A=matrix(c(1,-1,2,0,3,1), nrow=2)
1084 A
1085 #[,1] [,2] [,3]
1086 #[1,]    1    2    3
1087 #[2,]   -1    0    1
1088 #Perform SVD calculation
1089 msvd=svd(A)
1090 msvd
1091 msvd$d
1092 #[1] 3.784779 1.294390
1093 msvd$u
1094 #[,1]      [,2]
1095 #[1,] -0.9870875 -0.1601822
1096 #[2,] -0.1601822  0.9870875
1097 msvd$v
1098 #[,1]      [,2]
1099 #[1,] -0.2184817 -0.8863403
1100 #[2,] -0.5216090 -0.2475023
1101 #[3,] -0.8247362  0.3913356
1102 #One can verify that A=UDV', where V' is transpose of V.
1103 verim=msvd$u%*%diag(msvd$d)%*%t(msvd$v)
1104 verim
1105 #[,1]      [,2] [,3]
1106 #[1,] 1 2.000000e+00 3
1107 #[2,] -1 1.665335e-16 1
1108 round(verim)
1109 #[,1] [,2] [,3]
1110 #[1,] 1 2 3
1111 #[2,] -1 0 1
1112 #This is the original data matrix A
```

1113

---

1114 The covariance of the space-time anomaly matrix  $A$  is a spatial matrix:

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

1115 where  $Y$  is the number of columns of  $A$  and is equal to the length of time being  
 1116 considered.

1117 Strictly speaking, each column of an anomaly matrix  $A$  has a zero mean. But in  
 1118 practice, due to different normalization periods, the mean is not zero.

```
1119 covm = (1/(dim(A)[2]))*A%*%t(A)
1120 eigcov = eigen(covm)
1121 eigcov$values
1122 #[1] 4.7748518 0.5584816
1123 eigcov$vectors
1124 # [,1]      [,2]
1125 #[1,] -0.9870875 0.1601822
1126 #[2,] -0.1601822 -0.9870875
```

1127 Thus, the eigenvectors of a covariance matrix are the same as the SVD eigenvectors  
 1128 of the anomaly data matrix. The eigenvalues of the covariance matrix and the SVD  
 1129 have following relationship

```
1130 ((msvd$d)^2)/(dim(A)[2])
1131 #[1] 4.7748518 0.5584816
1132 eigcov$values
1133 #[1] 4.7748518 0.5584816
```

1134 Therefore, the EOFs from a given space-time dataset can be calculated directly by  
 1135 using an SVD command and do not need the step of calculating the covariance ma-  
 1136 trix. With efficient SVD algorithms, this shortcut can save significant amounts of  
 1137 time for an EOF analysis, also known as the principal component analysis (PCA) in  
 1138 the statistics community, compared to the traditional covariance matrix approach.  
 1139 Therefore, the result is extremely helpful for the EOF analysis, which is an indis-  
 1140 pensable modern tool of climate data analysis. The result is formally described as  
 1141 a theorem whose proof is also provided as follows.

1142 **Theorem 1.1** *The eigenvectors  $\mathbf{u}_k$  of the covariance matrix  $C = (1/Y)AA^t$*

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

1143 *are the same as the SVD spatial modes of  $A = UDV^t$ . The eigenvalues  $\lambda_k$  of  $C_{N \times N}$   
 1144 and the SVD eigenvalues  $d_k$  of  $A_{N \times Y}$  have the following relationship*

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

1145 *where  $Y$  is the total time length (i.e., time dimension) of the anomaly data matrix  
 1146  $A$ , and  $N$  is the total number of stations for  $A$  (i.e., space dimension).*

1147 **Proof** The SVD of the space-time data matrix is

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

<sup>1148</sup> The data matrix  $A$ 's corresponding covariance matrix is thus

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

<sup>1149</sup> In the above, we have used  $V^tV = I_Y$  is a  $Y \times Y$  identity matrix according to  
<sup>1150</sup> the SVD definition. The identity matrix's dimension is  $Y$  if  $Y \leq N$ . Otherwise,  
<sup>1151</sup>  $V^tV = I_N$ .

<sup>1152</sup> The expression  $C = \frac{1}{Y}UD^2U^t$  is the EOF expansion of the covariance matrix  
<sup>1153</sup> and means that a covariance matrix consists of EOFs and its associated variance,  
<sup>1154</sup> or "energy."

<sup>1155</sup> The covariance matrix's eigenvalue problem is

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

<sup>1156</sup> where

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

<sup>1157</sup> is the diagonal eigenvalue matrix with its elements as

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

<sup>1158</sup> In eq. (1.50), we used  $U^tU = I_Y$  based on the SVD definition. Equation (1.50)  
<sup>1159</sup> becomes an eigenvalues problem because  $\Lambda$  is a diagonal matrix:  $C\mathbf{u}_k = \lambda_k\mathbf{u}_k$  ( $k =$   
<sup>1160</sup>  $1, 2, \dots, Y$ ), where  $\mathbf{u}_k$  is the  $k$ th column vector of the space matrix  $U$ . Thus, eq.  
<sup>1161</sup> (1.50) implies the first part of the theorem: The space eigenvectors  $U$  from the  
<sup>1162</sup> SVD of the space-time data matrix  $A$  are the same as the eigenvectors of the  
<sup>1163</sup> corresponding covariance matrix  $C$ .

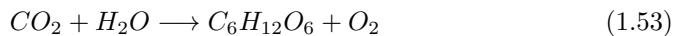
<sup>1164</sup> Equation (1.52) is exactly the second part of the theorem. The proof is thus  
<sup>1165</sup> complete.

<sup>1166</sup> In the above proof, we implicitly assumed that the columns of data matrix  $A$  are  
<sup>1167</sup> independent when  $Y \leq N$ . Independence of a set of vectors means that no one can  
<sup>1168</sup> be expressed in terms of the others, so no vectors can be replaced. Real climate  
<sup>1169</sup> data always satisfy this independence assumption.

## 1.7 Mass balance for chemical equations in marine chemistry

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

In the process of photosynthesis, plants convert the solar radiant energy carried by photons, plus carbon dioxide ( $\text{CO}_2$ ) and water ( $\text{H}_2\text{O}$ ), into glucose ( $\text{C}_6\text{H}_{12}\text{O}_6$ ) and oxygen ( $\text{O}_2$ ). The chemical equation for this conversion could be written schematically as



However, the conservation of mass requires that the atomic weights on both sides of the equation be equal. The photons have no mass. Thus, the chemical equation as written above is incorrect. The correct equation should specify precisely how many  $\text{CO}_2$  molecules react with how many  $\text{H}_2\text{O}$  molecules to generate how many  $\text{C}_6\text{H}_{12}\text{O}_6$  and  $\text{O}_2$  molecules. Suppose that these coefficients are  $x_1, x_2, x_3, x_4$ . We then have



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

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

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

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

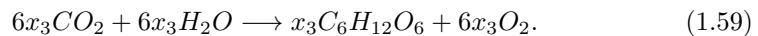
Similarly, the balance of oxygen atoms is

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

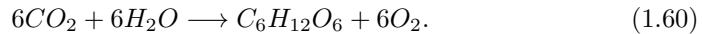
We thus have three equations in four variables. Thus, these equations have infinitely many solutions. We can set any variable fixed, and express the other three using this fixed variable. Since the largest molecule is glucose, we set its coefficient  $x_3$  fixed. Then we have

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

1199 Thus, the chemical equation is



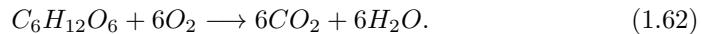
1200 If we want to produce one glucose molecule, i.e.,  $x_3 = 1$ , then we need 6 carbon  
1201 dioxide and 6 water molecules:



1202 Similarly, one can write chemical equations for many common reactions, such as  
1203 iron oxidation



1204 and the redox reaction in a human body which consumes glucose and converts the  
1205 glucose into energy, water and carbon dioxide



## 1206 1.8 Multivariate linear regression using matrix 1207 notations

---

1208

1209 This section is an application of matrix algebra in statistical data analysis, partic-  
1210 ularly on a linear regression for more than one variable. The linear regression in  
1211 Section ?? discussed the fitting of a straight line on the  $xy$ -plane  $y = b_1x + b_0$  to a  
1212 pair of data vectors of length  $N$ :  $[x_d]_{N \times 1}, [y_d]_{N \times 1}$ . It resulted in correlation, trend  
1213 and other regression quantities. An example is the correlation between the January  
1214 SOI as  $x$  and U.S. temperature as  $y$ . The non-trivial correlation can suggest that  
1215 there may be a physical mechanism to explain how the January SOI influences the  
1216 U.S. January temperature.

1217 Here, we use  $y = b_1x + b_0$  as the representation of a deterministic fitting func-  
1218 tion. This expression does not involve random variables. Most statistics books would  
1219 consider random linear models, which distinguishes between random variables and  
1220 their deterministic estimators by data. Our simple linear mathematical model for-  
1221 mulation here is equivalent to using only the deterministic estimators.

1222 The U.S. January temperature can be influenced by multiple factors in addition  
1223 to the SOI. These factors may include the North Atlantic sea surface temperature  
1224 (SST), the North Pacific SST, Arctic pressure conditions, etc. Then, the linear model  
1225 becomes

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

1226 When  $n = 1$ , this degenerates into the single variable regression.

1227 In the following we will present a few R examples of multivariate regression and  
1228 its applications. The mathematical theory behind the R code for a multivariate  
1229 regression deals mostly with the matrix operations, which can be found from any  
1230 standard textbooks covering multivariate regression.

1231

---

1232 **Example 1.9** This example shows a two-variable regression.

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

1233 Geometrically, this is an equation of a plane in a 3D space  $(x_1, x_2, y)$ . Given three  
 1234 points not on a straight line, a plane can be determined uniquely. This means specifying  
 1235 three  $x_1$  coordinate values, three  $x_2$  coordinate values and three  $y$  coordinate  
 1236 values, which can be done by means of the following R code:

```
1237 x1=c(1,2,3) #Given the coordinates of the 3 points
1238 x2=c(2,1,3)
1239 y=c(-1,2,1)
1240 df=data.frame(x1,x2,y) #Put data into the data.frame format
1241 fit <- lm(y ~ x1 + x2, data=df)
1242 fit#Show the regression results
1243 #Call:
1244 # lm(formula = y ~ x1 + x2, data = df)
1245 #Coefficients:
1246 # (Intercept)          x1            x2
1247 #-5.128e-16    1.667e+00   -1.333e+00
1248
1249 1.667*x1-1.333*x2 #Verify that 3 points determining a plane
1250 #[1] -0.999  2.001  1.002
```

1251 **Example 1.10** This example will show that four arbitrarily specified points can-  
 1252 not all be on a plane. The fitted plane has the shortest distance squares, i.e., the  
 1253 least squares (LS), or minimal mean square error (MMSE). Thus, the residuals are  
 1254 non-zero, in contrast to the zero residuals in the previous example.

```
1255 u=c(1,2,3,1)
1256 v=c(2,4,3,-1)
1257 w=c(1,-2,3,4)
1258 mydata=data.frame(u,v,w)
1259 myfit <- lm(w ~ u + v, data=mydata)
1260 summary(myfit)#Show the result
1261 #Coefficients:
1262 # Estimate Std. Error t value Pr(>|t|)
1263 #(Intercept) 1.0000 1.8708 0.535 0.687
1264 #u 2.0000 1.2472 1.604 0.355
1265 #v -1.5000 0.5528 -2.714 0.225
```

1266 **Example 1.11** This example will show a general multivariate linear regression  
 1267 using R. It has three independent variables, one dependent variable, and ten data  
 1268 points. For R program simplicity, the data are generated by an R random number  
 1269 generator. Again, R requires that the data be put into data frame format so that  
 1270 a user can clearly specify which are independent variables, also called explainable  
 1271 variables, and which is the dependent variable, also called response variable.

```
1272 dat=matrix(rnorm(40),nrow=10, dimnames=list(c(letters[1:10]), c(LETTERS[23:26])))
1273 fdat=data.frame(dat)
1274 fit=lm(Z~ W + X + Y, data=fdat)
```

```

1275 summary(fit)
1276
1277 #Coefficients\index{linear regression!coefficients}
1278 #             Estimate Std. Error t value Pr(>|t|)
1279 #(Intercept)   0.36680   0.16529   2.219   0.0683
1280 #W            0.11977   0.20782   0.576   0.5853
1281 #X            -0.53277  0.19378  -2.749   0.0333
1282 #Y            -0.04389  0.14601  -0.301   0.7739

```

1283 Thus, the linear model is

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

1284 The 95% confidence interval for  $W$ 's coefficient is  $0.12 \pm 2 \times 0.21$ , that for  $X$ 's  
 1285 coefficient is  $-0.53 \pm 2 \times 0.19$ ,  $Y$ 's coefficient is  $-0.04389 \pm 2 \times 0.15$ . Each confi-  
 1286 dence interval includes zero. Thus, there is no significant non-zero trend for the  $Z$   
 1287 data with respect to  $W, X, Y$ . This result is to be expected, because the data are  
 1288 randomly generated and thus should not have a trend.

1289 In practical applications, a user can simply convert the data into the same data  
 1290 frame format as shown here. Then, R command

```

1291 lm(formula = Z ~ W + X + Y, data = fdat)
1292 can do the regression job.
1293

```

---

1294 R can also do nonlinear regression with specified functions, such as quadratic  
 1295 functions and exponential functions. See examples from the URLs

```

1296 https://www.zoology.ubc.ca/~schluter/R/fit-model/
1297 https://stat.ethz.ch/R-manual/R-devel/library/stats/html/nls.html

```

## 1.9 Chapter summary

---

1298 This chapter introduces matrices and linear algebra from the perspective of space-  
 1299 time climate data, with rows corresponding to spatial locations (e.g., the order of  
 1300 grid boxes, grid point IDs, or climate station IDs) and columns corresponding to  
 1301 time (e.g., the month stamps January 1901, February 1901, ..., December 2017  
 1302 for the monthly data, or the day stamps 1 June 2017, 2 June 2017, ..., 30 June  
 1303 December 2017 for the daily data of June 2017). This way of introducing matrices  
 1304 as organized collections of data is different from that of most textbooks which  
 1305 describe a matrix as originating from a set of linear equations. This particular way  
 1306 to arrange a space-time climate data matrix follows the similar philosophy of the  
 1307 netCDF data file, which has recently become a very popular format for representing  
 1308 extensive climate datasets.

1311 The most essential methods of linear algebra have to do with matrix properties  
 1312 (e.g., eigenvalues, eigenvectors, and SVD matrices) and matrix operations (e.g.,

1313 solving linear equations, and making linear transformations). Our chapter has de-  
 1314 scribed the following linear algebra methods:

- 1315    (i) Matrix arithmetic: addition of a matrix plus another matrix, subtraction of  
 1316       a matrix minus another matrix, multiplication of a matrix times another  
 1317       matrix, and multiplication of a matrix by a scalar, and division of the  
 1318       identity matrix by a matrix which is a way to define the inverse matrix of  
 1319       the latter.
- 1320    (ii) An eigenvector  $w$  of a matrix  $C$  is such a special vector that  $Cw$  is in  $w$ 's  
 1321       own direction. Thus, there exists a scalar  $\lambda$  such that  $Cw = \lambda w$ , where  $\lambda$   
 1322       is called an eigenvalue of the matrix  $C$  corresponding to the eigenvector  $w$ .
- 1323    (iii) SVD decomposes the space-time data matrix  $A$  into three matrices: a dimen-  
 1324       sionless spatial orthogonal matrix  $U$  as spatial eigenvectors, a dimension-  
 1325       less temporal orthogonal matrix  $V$  as temporal eigenvectors, and a diagonal  
 1326       “energy level” matrix  $D$  related to the eigenvalues of  $A$ 's covariance matrix,  
 1327       i.e.,

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

1328 The column vectors of  $U$  may represent the intrinsic spatial patterns of  
 1329 the climate data (e.g., the warming pattern of the eastern tropical Pacific  
 1330 for the El Niño based on the SST data), and those of  $V$  may represent  
 1331 the intrinsic temporal patterns (e.g., the El Niño peaks in the time series  
 1332 determined by a column vector of  $V$ ). See the specific patterns shown in  
 1333 the examples of Section ?? on the SVD analysis for the reanalysis data.  
 1334 The diagonal values of  $D$  show the “energy level” (or the strength) of the  
 1335 corresponding spatial and temporal patterns. In the last few years, the SVD  
 1336 method has become more widely used in science and engineering because of  
 1337 its universal power in decomposing any space-time matrix. Mathematicians  
 1338 have now begun to modernize the proofs of the theorems of linear algebra  
 1339 by using the SVD approach, instead of the traditional approach known as  
 1340 echelon reduction through row operations.

- 1341    (iv) Mathematical models based on linear equations have been developed and  
 1342       employed for many climate science applications, such as the mass balance  
 1343       model for chemical equations.
- 1344    (v) Multivariate linear regression has been formulated here from the perspective  
 1345       of matrix algebra.
- 1346    (vi) R codes have been written for all these methods, which can be conveniently  
 1347       used for solving climate data analysis problems in research and practical  
 1348       applications.

1349 Although the basic materials of linear algebra appeared as early as the 17th  
 1350 century, widespread university classroom education in linear algebra in the United  
 1351 States began only in the 1950s and later (Tucker 1993). At that time, electronic  
 1352 computers began to become available to solve large numbers of linear equations for  
 1353 engineering applications. Today's vastly increased computing power now provides

1354 opportunities to modernize educational aspects of this branch of mathematics, such  
1355 as using SVD to compute various properties of a matrix and to prove numerous  
1356 theorems of linear algebra. Climate mathematics takes full advantage of this tech-  
1357 nical progress. Many graduate students in climate science now use EOFs and PCs  
1358 in their research, whereas the EOF-PC techniques were rarely taught as part of  
1359 climate science education before the 1990s.

## References and Further Readings

- <sup>1361</sup> [1] Golub, G.H. and C. Reinsch, 1970: Singular value decomposition and least squares solutions. *Numerische mathematik* 14, 403-420.

<sup>1362</sup>

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.

- <sup>1364</sup> [2] Lorenz, E.N., 1956: Empirical orthogonal functions and statistical weather prediction. *Scientific Report No. 1, Statistical Forecasting Project*. Air Force Research Laboratories, Office of Aerospace Research, USAF, Bedford, Massachusetts, 49pp.

<sup>1366</sup>

<sup>1367</sup>

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.

- <sup>1369</sup> [3] Strang, G., 2016: *Introduction to Linear Algebra*. 5th edition, Wellesley-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.

- <sup>1372</sup> [4] Tucker, A., 1993: The growing importance of linear algebra in undergraduate 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.

1374

1375

## Exercises

1376

**1.1** The following are the SVD results

```

1378      mat
1379          [,1] [,2]
1380      [1,]    1    1
1381      [2,]    1   -1
1382
1383      svd(mat)
1384      $d
1385      [1] 1.414214 1.414214
1386
1387      $u
1388          [,1]      [,2]
1389      [1,] -0.7071068 -0.7071068
1390      [2,] -0.7071068  0.7071068
1391
1392      $v
1393          [,1] [,2]
1394      [1,]    -1    0
1395      [2,]     0   -1

```

1396 Use  $A = UDV^t$  to recover the first column of

```

1397      mat
1398          [,1] [,2]
1399      [1,]    1    1
1400      [2,]    1   -1

```

1401 Show detailed calculations of all the relevant matrices and vectors, and use  
1402 space-time decomposition to describe your results.

1403 For extra credit: Describe the spatial and temporal modes, and their corre-  
1404 sponding variances or energies.

- 1405   **1.2** Use R and the updated Darwin and Tahiti standardized SLP data to repro-  
 1406       duce the EOFs and PCs and to plot the EOF pattern maps and PC time  
 1407       series.
- 1408   **1.3** Do the same procedures in the previous problem but for original non-standardized  
 1409       data. Comment on the difference of the results from the previous problem.
- 1410   **1.4** (a) Download the monthly precipitation data at five different stations over  
 1411       the United States from the USHCN website:
- 1412
- 1413       [http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn\\_map\\_interface.html](http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn_map_interface.html)
- 1414
- 1415       (b) Use R to organize the January data from 1951 to 2010 into the space-time  
 1416       format.
- 1417       (c) Compute the climatology of each station as the 1971-2000 mean.
- 1418       (d) Compute the space-time anomaly data matrix  $A$  as the original space-time  
 1419       data matrix minus the climatology.
- 1420       (e) Use R to make the SVD decomposition of the space-time anomaly data  
 1421       matrix  $A = UDV^t$ .
- 1422       (f) Output the  $U$  and  $D$  matrices.
- 1423   **1.5** In the previous problem, use R and the formula  $UDV^t$  to reconstruct the  
 1424       original data matrix  $A$ . This is a verification of the SVD decomposition, and  
 1425       is also called EOF-PC reconstruction.
- 1426   **1.6** Use R to plot the maps of the first three EOF modes from the  $U$  matrix in  
 1427       the previous problem in a way similar to the two El Niño mode maps shown  
 1428       in Fig. ???. Try to explain the climate meaning of the EOF maps.
- 1429   **1.7** Use R to plot the first three PC time series from the  $V$  matrix in Problems  
 1430       4.4 and 4.5. Try to explain the climate meaning of the time series.
- 1431   **1.8** (a) A covariance matrix  $C$  can be computed from a space-time observed  
 1432       anomaly data matrix  $X$  which has  $N$  rows for spatial locations and  $Y$  columns  
 1433       for time in years:

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

1434       This is an  $N \times N$  matrix. Choose an  $X$  data matrix from the USHCN  
 1435       annual total precipitation data at three California stations from north to  
 1436       south [Berkeley, CA (040693); Santa Barbara, CA (047902); Cuyamaca, CA  
 1437       (042239)] and five years from 2001 to 2005. Consider the anomaly data with  
 1438       respect to the 2001-2005 mean and use R to calculate a covariance matrix for  
 1439        $N = 3$  and  $Y = 5$ .

1440       (b) Use R to find the inverse matrix of the covariance matrix  $C$ .

1441       (c) Use R to find the eigenvalues and eigenvectors of  $C$ .

1442       (d) Use R to make SVD decomposition of the data matrix  $X = UDV^t$ . Ex-  
 1443       plicitly write out the three matrices  $U$ ,  $D$  and  $V$ .

1444       (e) Use R to explore the relationship between the eigenvalues of  $C$  and the  
 1445       matrix  $D$ .

- 1446 (f) Compare the eigenvectors of  $C$  and the matrix  $U$ .  
1447 (g) Plot the PC time series and describe their behavior.
- 1448  
1449 **1.9** The burning of methane ( $CH_4$ ) with oxygen ( $O_2$ ) produces water ( $H_2O$ ) and  
1450 carbon dioxide ( $CO_2$ ). Balance the chemical reaction equation.  
1451 **1.10** The burning of propane ( $C_3H_8$ ) with oxygen ( $O_2$ ) produces water ( $H_2O$ ) and  
1452 carbon dioxide ( $CO_2$ ). Balance the chemical reaction equation.  
1453 **1.11** The burning of gasoline ( $C_8H_{18}$ ) with oxygen ( $O_2$ ) produces water ( $H_2O$ )  
1454 and carbon dioxide ( $CO_2$ ). Balance the chemical reaction equation.

## 2

## Matrix Theory and Visualization

1457

1456 This chapter includes a slightly more advanced theory of matrix compared with the  
 1459 last. It includes (i) the concepts of independence, (ii) spanned spaces by multiple  
 1460 vectors, (iii) rank and other properties of a matrix, (iv) more on SVD, and (v)  
 1461 visualization of matrices and their decomposed vectors and singular values.

1462 Many mathematicians regard mathematics as “beautiful,” using terms such as  
 1463 “elegant,” “deep,” and “general.” Our book, however, pays more attention to *relevant*,  
 1464 *useful*, and *modern* (*RUM*) from the perspectives of both mathematical sci-  
 1465 ences and other fields. *Relevant* is reflected by that every matrix and its operations  
 1466 can have an interpretation story that can be easily understood by a layman. *Useful*  
 1467 is reflected by that each theory or method has non-trivial application examples in  
 1468 science or engineering. *Useful* is reflected in the extensive use of matrices in modern  
 1469 data science, machine learning, R or Python programming, which were not the case  
 1470 two or three decades ago. Matrix visualization is an example of modernness.

1471 This chapter continues to feature the space-time data arrangement, which uses  
 1472 rows of a matrix for spatial locations, and columns for temporal steps. This is the  
 1473 universal and fundamental information structure of our world. The singular value  
 1474 decomposition (SVD) helps reveal the spatial and temporal features of climate  
 1475 dynamics as singular vectors and the strength of their variability as singular values.

1476

1477

### 2.1 Matrix Definitions

1478 A matrix is a rectangular array of numbers (or even expressions), often denoted by  
 1479 an upper case letter in either boldface or plain **A** or *A*:

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

1480 This is an  $n \times p$  matrix, in which  $a_{ij}$  are called elements or entries of the matrix **A**,  $i$   
 1481 is the row index from 1 to  $n$  and  $j$  is the column index from 1 to  $p$ . The dimensions  
 1482 of a matrix are denoted by subscript, e.g.,  $\mathbf{A}_{n \times p}$  indicating that the matrix **A** has  $n$   
 1483 rows and  $p$  columns. The matrix may also be indicated by square brackets around

1484 a typical element  $\mathbf{A} = [a_{ij}]$  or sometimes  $\{A\}_{ij}$ , maybe even  $A_{ij}$ . If  $n = p$ , then the  
 1485 the array is a square matrix.

1486 Figure 2.1 is an example of a space-time climate data matrix. It is a subset of the  
 1487  $5^\circ \times 5^\circ$  gridded monthly surface air temperature anomalies from the NOAA Merged  
 1488 Land Ocean Global Surface Temperature Analysis (NOAAGlobalTemp) (Version  
 1489 4.0). The rows are indexed according to the spatial locations prescribed by the  
 1490 latitude and longitude of the centroid of a  $5^\circ \times 5^\circ$  grid box (See the entries of the  
 1491 first two columns in boldface). The columns are indexed according to time (See  
 1492 the first row entries in boldface). The other entries are the temperature anomalies  
 1493 with respect to the 1971-2000 monthly climatology. The anomalies are arranged  
 1494 according to the locations by rows and the time by columns. The units for the  
 1495 anomaly data are  $^\circ\text{C}$ .

1496 For a given month, the spatial temperature data on the Earth's surface is itself a  
 1497 2-dimensional array. To make a space-time matrix, we assign each grid box a unique  
 1498 index  $s$  from 1 to  $n$  if the spatial region has  $n$  grid boxes. The index assignment is  
 1499 subjective, depending on the application needs. The commonly used way is to fix a  
 1500 latitude and increase the index number as the longitude increases, as indicated by  
 1501 the first two columns of the data matrix shown in Figure 2.1. When the longitude  
 1502 is finished at this latitude band, go to the next latitude band until the completion  
 1503 of the latitude. This can go from south to north, from north to south. Of course,  
 1504 one can fix the longitude first, and increase the index according to the ascending  
 1505 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	<b>242.5</b>	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	<b>247.5</b>	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	<b>252.5</b>	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	<b>257.5</b>	-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	<b>262.5</b>	-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	<b>267.5</b>	-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	<b>272.5</b>	-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	<b>277.5</b>	-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	<b>282.5</b>	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	<b>287.5</b>	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	<b>292.5</b>	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	<b>297.5</b>	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	<b>302.5</b>	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	<b>307.5</b>	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	<b>312.5</b>	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	<b>317.5</b>	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	<b>322.5</b>	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	<b>327.5</b>	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.

1506 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  $36 \times 72$  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](http://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

1546

1547

1548 This section provides a concise list to summarize fundamental properties and the  
 1549 commonly used operations of matrices. We limit our materials to the basics that  
 1550 are sufficient for this book.

1551 (i) Zero matrix: A zero matrix has its every entry equal to zero:  $\mathbf{0} = [0]$ , or  $0 = [0]$ ,  
 1552 or explicitly

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

1553 (ii) Identity matrix: An identity matrix is a square matrix whose diagonal entries  
 1554 are all equal to one and whose off-diagonal entries are all equal to zero, and  
 1555 is denoted by  $I$  or  $\mathbf{I}$ . See an expression of an identity matrix below:

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

1556 People also use the following notation

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

1557 where

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

1558 is called the Kronecker delta.

1559 An identity matrix may be regarded as a special case of a *diagonal matrix*,  
 1560 which refers to any square matrix whose off-diagonal elements are zero.  
 1561 Hence, a diagonal matrix has the following general expression:

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

1562 where  $d_1, d_2, \dots, d_n$  are the n-diagonal elements. If  $d_i = 1, i = 1, 2, \dots, n$ ,  
 1563 then the diagonal matrix becomes an identity matrix.

1564 (iii) A transpose matrix: The *transpose* of a matrix  $A$  is obtained by interchanging  
 1565 the rows and columns. The new matrix is denoted by  $\mathbf{A}^t$ . Computing the  
 1566 matrix transpose is very easy, simply rotating each horizontal row clockwise  
 1567 to a vertical column, one row at a time. If  $\mathbf{A}$  has dimension  $n \times p$ , its

1568 transpose has dimension  $p \times n$ . The elements of the transposed matrix are  
 1569 related to the originals by

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

1570 For example, if

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

1571 then

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

1572 The equation  $\mathbf{A}^t = \mathbf{A}$  or  $a_{ij} = a_{ji}$  is equivalent to that the matrix  $\mathbf{A}$  is  
 1573 symmetric. Of course, a symmetric matrix must be a square matrix.

- 1574 (iv) Equal matrices: Two matrices  $\mathbf{A}$  and  $\mathbf{B}$  are equal if every pair of corresponding  
 1575 entries are equal, i.e., the equation  $\mathbf{A} = \mathbf{B}$  is equivalent to  $a_{ij} = b_{ij}$  for all  
 1576  $i$  and  $j$ .
- 1577 (v) Matrix addition: The sum of two matrices  $\mathbf{A}$  and  $\mathbf{B}$  is defined by the sum of  
 1578 corresponding entries, i.e.,  $\mathbf{A} + \mathbf{B} = [a_{ij} + b_{ij}]$ .
- 1579 (vi) Matrix subtraction: The difference of two matrices  $\mathbf{A}$  and  $\mathbf{B}$  is defined by the  
 1580 difference of corresponding entries, i.e.,  $\mathbf{A} - \mathbf{B} = [a_{ij} - b_{ij}]$ . The equation  
 1581 of  $\mathbf{A} = \mathbf{B}$  is equivalent to  $\mathbf{A} - \mathbf{B} = 0$ .
- 1582 (vii) Row vector: A row vector is of dimension  $p$  is a  $1 \times p$  matrix:

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

1583 Matrix  $\mathbf{A}_{n \times p}$  may be regarded as a stack of  $n$  row vectors  $\mathbf{a}_{i:}, i = 1, 2, \dots, n$ ,  
 1584 each of which is a  $p$ -dimensional row vector. Hence,

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

1585 Here, : in the second position of the double-index subscript means the  
 1586 inclusion of all the columns.

- 1587 (viii) Column vector: A column vector of dimension  $n$  is an  $n \times 1$  matrix

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

1588 The transpose of a column vector becomes a row vector, and vice versa.

1589 Matrix  $\mathbf{A}_{n \times p}$  may be regarded as an array of  $p$  column vectors  $\mathbf{a}_{:,j}, j =$   
 1590  $1, 2, \dots, p$ , each of which is an  $n$ -dimensional column vector. Hence,

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

1591 Here, : in the first position of the double-index subscript means the inclusion  
 1592 of all the rows.

1593 (ix) Dot product of two vectors: Two vectors of the same dimension can form a  
 1594 *dot product* which is equal to the sum of the products of the corresponding  
 1595 entries:

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

1596 The dot product is also called an *inner product*.

1597 As example, if

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

1598 then

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

1599 The *amplitude* of vector  $\mathbf{u}$  of dimension  $n$  is defined as

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

1600 Sometimes, the amplitude is also called length, or Euclidean length, or  
 1601 magnitude. Please do not mix the concept of Euclidean length of a vector  
 1602 with the dimensional length of a vector. The latter means the number of  
 1603 entries of a vector, i.e.,  $n$ .

1604 If the Euclidean length of  $\mathbf{u}$  is equal to one, we say that the  $\mathbf{u}$  is a *unit*  
 1605 *vector*. If every element of  $\mathbf{u}$  is zero, then we say that  $\mathbf{u}$  is a *zero vector*.

1606 By the definition of dot product, we have

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

1607 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  
 1608  $|\mathbf{u}| = |\mathbf{v}| = 1$ , then we say that  $\mathbf{u}$  and  $\mathbf{v}$  are *orthonormal*.

1609 (x) Matrix multiplication: The product of matrix  $\mathbf{A}_{n \times p}$  and matrix  $\mathbf{B}_{p \times m}$  is an  
 1610  $n \times m$  matrix  $\mathbf{C}_{n \times m}$  whose element  $c_{ij}$  is the dot product of the  $i$ th row  
 1611 vector of  $A$  and  $j$ th column vector of  $B$ :

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

1612 We denote

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

1613 or

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

1614 Note that the number of columns of  $\mathbf{A}$  and that of rows of  $\mathbf{B}$  must be the  
 1615 same when the multiplication  $\mathbf{AB}$  can be made, because the dot product  
 1616  $\mathbf{a}_{i:} \cdot \mathbf{b}_{:j}$  requires this condition. This is referred to as the dimension-matching  
 1617 condition for matrix multiplication. If this condition is violated, the two

1618 matrices cannot be multiplied. For example, for the following two matrices

1619

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

1620

We can compute

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

1621

However, the expression

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

1622

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

1626

1627

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

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

1628

1629

1630

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

1631

1632

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

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

1633

1634

1635

The scalar multiplication can be extended to multiple vectors

$(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p)$

or matrices to form a *linear combination*:

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

1636

1637

1638

1639

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

1640

1641

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

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

1642

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

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

1643 Not all the matrices have an inverse. If a matrix has an inverse, then the  
1644 matrix is said to be *invertible*. Equivalently,  $\mathbf{A}^{-1}$  exists.

1645 As an example of the matrix inversion, given

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

1646 we have

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

1647 Hand-calculation for the inverse of a small matrix is already very difficult,  
1648 and that for the inverse of a large matrix is almost impossible. Computers  
1649 can do the calculations for us, as shown in examples later.

1650 According to the definition of inverse, we have the following formula for  
1651 the inverse of the product of two matrices:

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

1652 if both  $\mathbf{A}$  and  $\mathbf{B}$  are invertible matrices. Please note the order switch of  
1653 the matrices.

1654 With the definition of an inverse matrix, we can define the *matrix division*  
1655 by

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

1656 when  $\mathbf{B}^{-1}$  exists. In matrix operations, we usually do not use the concept  
1657 of matrix division, but always use the matrix inverse and matrix multipli-  
1658 cation.

1659 (xii) More properties of the matrix transpose:

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

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

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

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

1660 (xiii) Orthogonal matrices: An *orthogonal matrix*<sup>1</sup> is one whose row vectors are  
1661 orthonormal. In this case, the inverse matrix can be easily found: It is its  
1662 transpose. That is, if  $\mathbf{A}$  is an orthogonal matrix, then

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

1663 The proof of this claim is very simple. The orthonormal property of the  
1664 row vectors of  $\mathbf{A}$  implies that

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

1665 By the definition of matrix inverse,  $\mathbf{A}^t$  is the inverse matrix of  $\mathbf{A}$ .

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

1666 If  $\mathbf{A}$  is an orthogonal matrix, its row vectors are also orthonormal. This  
 1667 can be proved by multiplying both sides of the above by  $\mathbf{A}^t$  from the left:

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

1668 Then multiply both sides of this equation by  $(\mathbf{A}^t)^{-1}$  from the right:

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

1669 which yields

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

1670 This implies that the column vectors of  $\mathbf{A}$  are orthonormal.

1671 As an example, the following matrix

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

1672 is an orthogonal matrix for any given real number  $\theta$ . You can easily verify  
 1673 this using the trigonometrical identity  $\sin^2 \theta + \cos^2 \theta = 1$ . We thus have

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

1674 You can easily verify that  $\mathbf{T}^{-1} \mathbf{T} = \mathbf{I}$  by hand-calculation of the product  
 1675 of the two matrices in this equation.

## 1676 2.3 Some basic concepts and theories of linear algebra

---

1677 According to Encyclopedia.com, “linear algebra originated as the study of linear  
 1678 equations.” Linear algebra deals with vectors, matrices and vector spaces. Before  
 1679 the 1950s, it was part of Abstract Algebra (Tucker 1993). In 1965 the Commit-  
 1680 tee on the Undergraduate Program in Mathematics, Mathematical Association of  
 1681 America, outlined the following topics for a stand-alone linear algebra course: Lin-  
 1682 ear systems, matrices, vectors, linear transformations, unitary geometry with char-  
 1683 acteristic values. The vectors and matrices have been dealt in the previous two  
 1684 sections of this chapter. This section deals with linear systems of equations, and  
 1685 linear transformations, and next with characteristic values.

### 1687 2.3.1 Linear equations

---

1688 A meteorologist needs to make a decision on what instruments to order under the  
 1689 following constraint. She is given a budget of \$10,000 to purchase 30 instruments for  
 1690 her observational sites. Her supplier has two products for the instrument: The first  
 1691 is \$30 per set, and the second \$40 per set. She would like to buy the second type of  
 1692 instrument as many as possible under the budget constraint. Then, the question is

1694 how many instruments of the second kind she can buy? This problem leads to the  
 1695 following linear system of two equations:

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

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

1696 The solution to these linear equations is  $x_1 = 20$  and  $x_2 = 10$ .

1697 This system of linear equations can be expressed in matrix and vectors as follows:

1698

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

1699 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)

1700 Then, the solution of this system may be tightly expressed in the following way:

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

1701 This expression is convenient for mathematical proofs, but is rarely used for solving  
 1702 a linear system, because finding an inverse matrix is computationally costly. A way  
 1703 to solve a linear system is to use Gauss elimination. The corresponding computing  
 1704 procedure is called the row operation on a matrix. There are numerous ways of  
 1705 solving a linear system. Some are particularly efficient for a certain system, such  
 1706 as a sparse matrix or a matrix of a diagonal bend of width equal 3 or 5. Efficient  
 1707 algorithms for a linear system, particularly an extremely large system, are forever  
 1708 a research topic. In this book we use a computer to solve a linear system without  
 1709 studying the algorithm details. The R and Python commands are below:

```
1710 solve(A, b) #R code for finding x
1711 numpy.linalg.solve(A, b) #Python code
```

### 2.3.2 Linear transformations

---

1712  
 1713 A *linear transformation* is to convert vector  $\mathbf{x}_{n \times 1}$  into  $\mathbf{y}_{m \times 1}$  using the multiplication  
 1714 of a matrix  $\mathbf{T}_{m \times n}$ :

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

1715 For example, the matrix

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

<sup>1717</sup> transforms the vector

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

<sup>1718</sup> into

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

<sup>1719</sup> This is the solution of the linear system in the previous subsection.

<sup>1720</sup> Usually, the linear transformation  $\mathbf{T}\mathbf{x}$  changes both direction and magnitude of  
<sup>1721</sup>  $\mathbf{x}$ . However, if  $\mathbf{T}$  is an orthogonal matrix, then  $\mathbf{T}\mathbf{x}$  does not change the magnitude of  
<sup>1722</sup>  $\mathbf{x}$ , and changes only the direction. This claim can be simply proved by the following  
<sup>1723</sup> formula:

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

<sup>1724</sup> Thus, if  $\mathbf{T}$  is an orthogonal matrix, then  $\mathbf{T}\mathbf{x}$  is a rotation of the vector  $\mathbf{x}$ .

### <sup>1725</sup> 2.3.3 Linear independence

---

<sup>1726</sup> The vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$  are *linearly dependent* if no vector can be represented  
<sup>1727</sup> by a linear combination of other  $p - 1$  vectors in this group. Otherwise, the group  
<sup>1728</sup> of vectors are linearly independent.

<sup>1729</sup> If it is not linearly independent, then there must be a vector which can be rep-  
<sup>1730</sup> resented by the other vectors through a linear combination. Suppose this vector is  
<sup>1731</sup>  $\mathbf{x}_1$ , then

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

<sup>1732</sup> where at least one of the coefficients  $d_2, d_3, \dots, d_p$  is non-zero. Thus, the linear  
<sup>1733</sup> system of equations for  $c_1, c_2, c_3, \dots, c_p$

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

<sup>1734</sup> has a non-zero solution. This system can be written as a matrix form

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

<sup>1735</sup> where column vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$  form the matrix  $\mathbf{X}$

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

<sup>1736</sup> the unknown vector is  $\mathbf{c}$

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

1738 and  $\mathbf{0}$  is the p-dimensional zero column vector. The solution of this matrix equation  
 1739 is

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

1740 However,  $\mathbf{c}$  must not be zero. This contradiction implies that  $\mathbf{X}^{-1}$  does not exist if  
 1741 the the column vectors are linearly dependent. In other words, if  $\mathbf{X}^{-1}$  exists, then  
 1742 its column vectors are linearly independent.

1743 Consider vectors in a three dimensional space. Any two column vectors  $\mathbf{x}_2$  and  
 1744  $\mathbf{x}_3$  define a plane. If  $\mathbf{x}_1$  can be written as a linear combination of  $\mathbf{x}_2$  and  $\mathbf{x}_3$ , then  
 1745 it must lie in the same plane. So,  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  and  $\mathbf{x}_3$  are linearly dependent. The matrix  
 1746  $[\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_3]_{3 \times 3}$  is not invertible.

### 1747 2.3.4 Determinants

---

1749 For a square matrix  $\mathbf{A}$ , a convenient notation and concept is its *determinant*. It is  
 1750 a scalar and is denoted by  $\det(\mathbf{A})$  or  $|\mathbf{A}|$ . For a  $2 \times 2$  matrix

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

1751 its determinant is

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

1752 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  
 1753 matrix, say  $\mathbf{A}_{172 \times 172}$ . The computer command for computing the determinant of a  
 1754 small square matrix is as follows:

```
1756 det(A) #R command for determinant
1757 np.linalg.det(a) #Python command for determinant
```

1758 Two 2-dimensional column vectors  $\mathbf{x}_1$  and  $\mathbf{x}_2$  can span a parallelogram, whose  
 1759 area  $S$  is equal to the absolute value of the determinant of the matrix consisting of  
 1760 the two vectors  $\mathbf{A} = [\mathbf{x}_1 \ \mathbf{x}_2]$ :

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

1761 Three 3-dimensional column vectors  $\mathbf{x}_1, \mathbf{x}_2$  and  $\mathbf{x}_3$  can span a parallelepiped,  
 1762 whose volume  $V$  is equal to the absolute value of the determinant of the matrix  
 1763 consisting of the three vectors  $\mathbf{A} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_3]$ :

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

1764 A few commonly used properties of determinant are listed below

- 1765 (a) The determinant of a diagonal matrix is the product of its diagonal elements.
- 1766 (b) If a determinant has a zero row or column, the determinant is zero.
- 1767 (c) The determinant does not change after a matrix transpose, i.e.,  $\det[\mathbf{A}^t] =$   
 1768  $\det[\mathbf{A}]$ .

- 1769 (d) The determinant of the product of two matrices:  $\det[\mathbf{AB}] = \det[\mathbf{A}]\det[\mathbf{B}]$ .  
 1770 (e) The determinant of the product of a matrix with a scalar:  $\det[c\mathbf{B}] = c^n\det[\mathbf{B}]$ ,  
 1771 if  $\mathbf{A}$  is an  $n \times n$  matrix.  
 1772 (f) The determinant of an orthogonal matrix is equal to 1 or -1.

### 2.3.5 Rank of a matrix

1775 The *rank* of  $\mathbf{A}$  is the greatest number of columns of the matrix that are linearly  
 1776 independent, and is denoted by  $r[\mathbf{A}]$ .

1777 For a  $3 \times 3$  matrix  $\mathbf{A}$ , if we treat each column as a 3-dimensional vector. If all  
 1778 three lie along a line (i.e., collinear), the rank of  $\mathbf{A}$  is one. If all three of the vectors  
 1779 lie in a plane, but are not collinear, the rank is two. If none of the three are collinear  
 1780 or lie in a plane, the rank is three.

1781 If the rank of a square matrix is less than its dimension, then at least one column  
 1782 can be a linear combinations of other columns, which implies that the determinant  
 1783 vanishes. If  $\mathbf{A}$  has rank  $r$ , it is possible to find  $r$  linearly independent columns, and  
 1784 all the other columns are linear combinations of these  $r$  independent columns.

1785 Some properties about the matrix rank are listed below:

- 1786 (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  
 1787 said to be *nonsingular*.  
 1788 (b) If  $\det[\mathbf{A}] = 0$ , then the rank of  $\mathbf{A}$  is less than  $n$ , and  $\mathbf{A}$  is not invertible and  
 1789 is said to be *singular*.  
 1790 (c) If  $\mathbf{B}$  is multiplied by a nonsingular matrix  $\mathbf{A}$ , the product has the same rank  
 1791 as  $\mathbf{B}$ .  
 1792 (d)  $0 \leq r[\mathbf{A}_{n \times p}] \leq \min(n, p)$ .  
 1793 (e)  $r[\mathbf{A}] = r[\mathbf{A}^t]$ .  
 1794 (f)  $r[\mathbf{AB}] \leq \min(r[\mathbf{A}], r[\mathbf{B}])$ .  
 1795 (g)  $r[\mathbf{AA}^t] = r[\mathbf{A}^t \mathbf{A}] = r[\mathbf{A}]$ .  
 1796 (h)  $r[\mathbf{A} + \mathbf{B}] \leq r[\mathbf{A}] + r[\mathbf{B}]$ .

1797 Computers can easily demonstrate the matrix computations following the theories  
 1798 presented in this chapter so far. The computer code is below.

```
1799 #R code: Computational examples of matrices
1800 A = matrix(c(1,0,0,4,3, 2), nrow = 3, byrow = TRUE)
1801 B = matrix(c(0,1,-1,2), nrow = 2) #form a matrix by columns
1802 C = A%*%B #matrix multiplication
1803 C
1804 #[1,]    0   -1
1805 #[2,]    4    8
1806 #[3,]    2    1
1807 t(C) # transpose matrix of C
1808 #[1,]    0    4    2
1809 #[2,]   -1    8    1
1810
1811 A = matrix(c(1, -1, 1, 2), nrow = 2, byrow = TRUE)
1812 solve(A) #compute the inverse of A
1813 #[1,] 0.6666667 0.3333333
```

```
1814 #[2,] -0.3333333 0.3333333
1815 A%*%solve(A) #verify the inverse of A
1816 #[1,] 1.000000e+00 0
1817 #[2,] 1.110223e-16 1
1818
1819 #Solve linear equations
1820 A = matrix(c(30, 40, 1, 1), nrow = 2, byrow = TRUE)
1821 b = c(1000, 30)
1822 solve(A,b)
1823 #[1] 20 10
1824 solve(A)%*%b #Another way to solve the equations
1825 det(A) #compute the determinant
1826 #[1] -10
1827
1828 library(Matrix)
1829 rankMatrix(A) #Find the rank of a matrix
1830 #[1] 2 #rank(A) = 2
1831
1832 #Orthogonal matrices
1833 p = sqrt(2)/2
1834 Q = matrix(c(p,-p,p,p), nrow=2)
1835 Q #is an orthogonal matrix
1836 # [,1] [,2]
1837 #[1,] 0.7071068 0.7071068
1838 #[2,] -0.7071068 0.7071068
1839 Q%*%t(Q) #verify 0 as an orthogonal matrix
1840 # [,1] [,2]
1841 #[1,] 1 0
1842 #[2,] 0 1
1843 det(Q) #The determinant of an orthogonal matrix is 1 or -1
1844 #[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

```

1846

## 2.4 Eigenvectors and eigenvalues

1847

1848  
1849

### 2.4.1 Definition of eigenvectors and eigenvalues

1850 The linear transform  $\mathbf{y} = \mathbf{Au}$  usually results in  $\mathbf{y}$  not parallel to  $\mathbf{u}$ . For example,

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

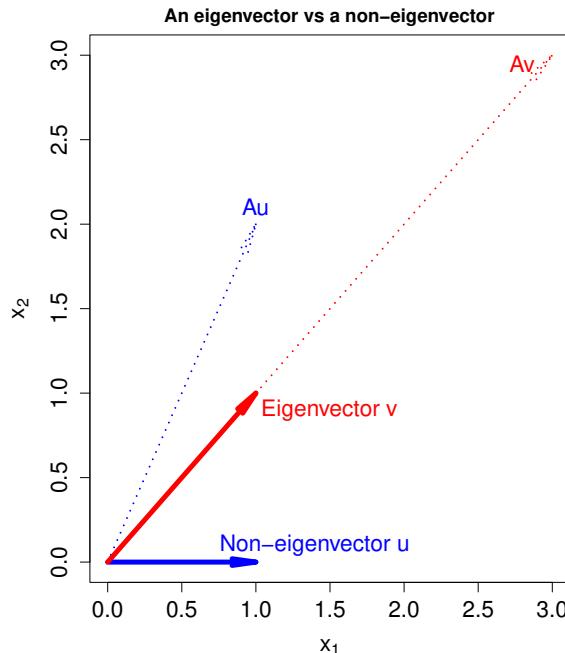
1851 The vectors  $\mathbf{u} = (1, 0)$  and  $\mathbf{Au} = (1, 2)$  are not parallel in the 2-dimensional space.  
1852 See the blue vectors in Fig. 2.2.

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.

1853 However, there exist some special vectors  $\mathbf{v}$  such that  $\mathbf{Av}$  is parallel to  $\mathbf{v}$ .  
1854 For example,  $\mathbf{v} = (1, 1)$  is such a vector, since  $\mathbf{Av} = (3, 3)$  is in the same direction  
1855 as  $\mathbf{v} = (1, 1)$ . See the red vectors in Fig. 2.2. If two vectors are parallel, then one  
1856 vector is a scalar multiplication of the other, e.g.,  $(3, 3) = 3(1, 1)$ . We denote this  
1857 scalar by  $\lambda$ . Thus,

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

1858 These vectors  $v$  are special to  $A$ , maintain their own orientation when multiplied  
1859 by  $A$ , and are called *eigenvectors*. Here, “eigen” is from German, meaning “self”,

<sup>1860</sup> “own”, “particular”, or “special.” <sup>2</sup> The corresponding scalars  $\lambda$  are called *eigenvalues*.  
<sup>1861</sup> The formula (2.69) is a mathematical definition of the eigenvalue problem  
<sup>1862</sup> for matrix  $\mathbf{A}$ .

<sup>1863</sup> If  $\mathbf{v}$  is an eigenvector, then its multiplication to a scalar  $c$  is also an eigenvector,  
<sup>1864</sup> since

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

<sup>1865</sup> Namely, all the vectors in the same direction  $\mathbf{v}$  are also eigenvectors. The eigenvectors  
<sup>1866</sup> of length equal one are called unit eigenvectors, or unitary eigenvectors, and  
<sup>1867</sup> are unique up to a positive or negative sign. Most computer programs output unit  
<sup>1868</sup> eigenvectors. Thus, eigenvector describes a direction or an orientation. Each square  
<sup>1869</sup> matrix has its own special orientations.

<sup>1870</sup> The aforementioned vector

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

<sup>1871</sup> is en eigenvector that maintains its own direction after multiplied by  $\mathbf{A}$ :

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

<sup>1872</sup> Thus,

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

<sup>1873</sup> is an eigenvector of  $A$  and  $\lambda = 3$  is an eigenvalue of  $A$ . The corresponding unit  
<sup>1874</sup> eigenvector is

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

<sup>1875</sup> Another eigenvector for the above matrix  $\mathbf{A}$  is  $\mathbf{v}_2 = (1, -1)$ :

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

<sup>1876</sup> The second eigenvalue is  $\lambda_2 = -1$ .

<sup>1877</sup> The computer code for the eigenvalues and eigenvectors of the above matrix  $\mathbf{A}$   
<sup>1878</sup> is below.

```
1879 #R code for eigenvectors and eigenvalues
1880 A = matrix(c(1, 2, 2, 1), nrow=2)
1881 eigen(A)
```

<sup>2</sup> 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.

```

1882  ##$values
1883  #[1] 3 -1
1884  ##$vectors
1885  #[,1]      [,2]
1886  #[1,] 0.7071068 -0.7071068
1887  #[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]]))

```

1888  
1889     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

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

1890     The non-zero solution  $\mathbf{v}$  of this equation requires that

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

1891     Expanding the determinant out leads to an  $N$ th degree polynomial in  $\lambda$ , which has  
1892     exactly  $N$  roots. Some roots may be repeated and hence counted multiple times  
1893     toward  $N$ . Some roots may be complex numbers, which are not discussed in this  
1894     book. Each root is an eigenvalue and corresponds to an eigenvector.

1895     The above concise determinant expression is tidy and useful for mathematical  
1896     proofs, but is not used as a computer algorithm for calculating eigenvalues or eigen-  
1897     vectors, because it is computationally costly or even impossible for a large matrix.  
1898     Nonetheless, some traditional textbooks of linear algebra defined eigenvalues using  
1899     Eq. (2.74). This traditional definition of eigenvalues appears to be abstract. The  
1900     eigenvector equation (2.73) seems not providing an image of the eigenvector. Thus,  
1901     many students forget the definition of eigenvalues and eigenvectors shortly after the  
1902     final exams of a linear algebra course.

1903     Figure 2.2 may be generated by the following computer code.

```

1904
1905 #R plot eigenvector v vs a non-eigenvector u
1906 #Create your working directory named LinAlg and go there
1907 setwd('/Users/ssten/LinAlg')
1908 getwd() #Verify that you are in the directory/folder
1909 #[1] "/Users/ssten/LinAlg"
1910
1911 setEPS() #Plot the figure and save the file
1912 postscript("fig0502.eps", width = 6)
1913 par(mar=c(4.5,4.5,2.0,0.5))
1914 plot(9,9,
1915       main = 'An_eigenvector_vs_a_non-eigenvector',
1916       cex.axis = 1.4, cex.lab = 1.4,
1917       xlim = c(0,3), ylim=c(0,3),
1918       xlab = bquote(x[1]), ylab = bquote(x[2]))

```

```

1919 arrows(0,0, 1,0, length = 0.25,
1920           angle = 8, lwd = 5, col = 'blue')
1921 arrows(0,0, 1,2, length = 0.3,
1922           angle = 8, lwd = 2, col = 'blue', lty = 3)
1923 arrows(0,0, 1,1, length = 0.25,
1924           angle = 8, lwd = 5, col='red')
1925 arrows(0,0, 3,3, length = 0.3,
1926           angle = 8, lwd = 2, col='red', lty = 3)
1927 text(1.4,0.1, 'Non-eigenvectoru', cex =1.4, col = 'blue')
1928 text(1.0,2.1, 'Au', cex =1.4, col = 'blue')
1929 text(1.5,0.9, 'Eigenvectorv', cex =1.4, col = 'red')
1930 text(2.8, 2.95, 'Av', cex =1.4, col = 'red')
1931 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('Aneigenvectorvvsanon-eigenvectoru')
plt.xlabel(r'$x_1$', fontsize = 25)
plt.ylabel(r'$x_2$', fontsize = 25)
plt.text(0.6, 0.1, 'Non-eigenvectoru',
         color = 'blue', fontsize =25)
plt.text(0.8, 2.0, 'Au',
         color = 'blue', fontsize =25)
plt.text(1.03,0.85, 'Eigenvectorv',
         color = 'red', fontsize =25)
plt.text(2.7, 2.9, 'Av',
         color = 'red', fontsize =25)
plt.show()

```

---

### 1933 2.4.2 Properties of eigenvectors and eigenvalues for a symmetric matrix

---

1936 A covariance matrix or a correlation matrix is a symmetric matrix and is often used  
 1937 in climate science. For a symmetric matrix  $\mathbf{A}$ , its eigenvalues and eigenvectors have  
 1938 the following properties:

- 1939 (a) Eigenvalues of a symmetric matrix are real numbers.
- 1940 (b) The  $n$  different unit eigenvectors of a symmetric matrix  $\mathbf{A}_{n \times n}$  are independent  
 1941 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  
 1942 any two different unit eigenvectors. We can use these unit vectors to express  
 1943 any  $n$ -dimensional vector using a linear combination:

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

1944 where  $c_1, c_2, \dots, c_n$  are scalar coefficients of the linear combination.

- 1945 (c) If all the eigenvalues of  $\mathbf{A}_{n \times n}$  are positive, then the *quadratic form*

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

1946 is also a positive scalar for any non-zero vector  $\mathbf{x}$ . A quadratic form may  
 1947 be used to express kinetic energy or total variance in climate science. The  
 1948 kinetic energy in climate science is always positive. For all the unit vectors  
 1949  $\mathbf{x}$ , the maximum of the quadratic form is equal to the largest eigenvalue of  
 1950  $\mathbf{A}$ . The maximum is achieved when  $\mathbf{x}$  is the corresponding eigenvector of  
 1951  $\mathbf{A}$ .

- 1952 (d) The rank of matrix  $\mathbf{A}$  is equal to the number of nonzero eigenvalues, where  
 1953 the multiplicity of repeated eigenvalues is counted.
- 1954 (e) Eigenvalues of a diagonal matrix are equal to the diagonal elements.
- 1955 (f) For a symmetric matrix  $\mathbf{A}_{n \times n}$ , its  $n$  unit column eigenvectors form an or-  
 1956 thogonal matrix  $\mathbf{Q}_{n \times n}$  such that  $\mathbf{A}_{n \times n}$  can be diagonalized by  $\mathbf{Q}_{n \times n}$  in  
 1957 the following way

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

1958 where  $\mathbf{D}$  is a diagonal matrix whose diagonal elements are eigenvalues  
 1959  $\lambda_1, \lambda_2, \dots, \lambda_n$  of  $\mathbf{A}$ . Further, the column vectors of  $\mathbf{Q}$  are the unit eigen-  
 1960 vectors of  $\mathbf{A}$ . This is a matrix diagonalization process.

1961 For the symmetric matrix  $\mathbf{A}$  in Eq. (4.51)

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

1962 its eigenvalues and eigenvectors are

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

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

<sup>1963</sup> Thus, the orthogonal matrix  $\mathbf{A}$  and the diagonal matrix  $\mathbf{d}$  are as follows

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

<sup>1964</sup> With  $\mathbf{Q}$ ,  $\mathbf{A}$ , and  $\mathbf{D}$ , you can easily verify Eq. (2.77) by hand calculation or by <sup>1965</sup> computer coding.

<sup>1966</sup> Equation (2.77) can be written in the following way

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

<sup>1967</sup> or

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

<sup>1968</sup> This is a process of matrix decomposition by orthogonal matrices or by eigenvectors.

<sup>1969</sup> Further, if all the eigenvalues are non-negative, then the matrix is said to be *positive* <sup>1970</sup> *semi-definite*, and the above formula can be written as

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

<sup>1971</sup> where  $\mathbf{v}^{(k)} = \sqrt{\lambda_k} \mathbf{q}^{(k)}$  ( $k = 1, 2, \dots, n$ ). The sample covariance matrix in climate <sup>1972</sup> science satisfies the positive eigenvalue assumption, and will be discussed in more <sup>1973</sup> details in the next chapter. The last equation means that a positive semi-definite <sup>1974</sup> symmetric matrix  $\mathbf{A}_{n \times n}$  can be decomposed into a sum of  $n$  outer products of <sup>1975</sup> eigenvectors.

<sup>1976</sup> Matrix  $\mathbf{A}$

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

<sup>1977</sup> is not positive semi-definite since its second eigenvalue is -1, but matrix  $\mathbf{C}$

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

<sup>1978</sup> is positive semi-definite since its eigenvalues are 3 and 1, which are positive. The <sup>1979</sup> matrix  $\mathbf{C}$  is actually positive definite. The eigenvectors are the same those of  $\mathbf{A}$ . <sup>1980</sup> Thus,

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

<sup>1981</sup> where  $\mathbf{Q}$  is given by Eq. (2.81), and  $\mathbf{D}_c$  is

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

<sup>1982</sup> The above can be verified by the following computer code.

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

```

1986  #$values
1987  #[1] 3 1
1988  #$vectors
1989  # [,1]      [,2]
1990  #[1,] 0.7071068 -0.7071068
1991  #[2,] 0.7071068  0.7071068
1992  Q = eigen(C)$vectors
1993  D = t(Q) %*% C %*% Q #Matrix diagonalization
1994  D
1995  #[1,]    3     0
1996  #[2,]    0     1
1997  Q %*% D %*% t(Q) #Matrix decomposition
1998  #[1,]    2     1
1999  #[2,]    1     2
2000  D[1,1]*Q[,1] %*% t(Q[,1]) + D[2,2]*Q[,2] %*% t(Q[,2])
2001  #[1,]    2     1
2002  #[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 = [[ 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

```

2003

2004

## 2.5 Hadamard and Other Matrix Multiplications

2005

- 2006 In addition to the regular matrix products and sweeping by multiplication described  
 2007 in Chapter 1, there still other matrix products, such as Hadamard product of two  
 2008 matrices. This section presents a few more matrix operations.

2009                   **2.5.1 Hadamard Product of Two Matrices of the Same**  
 2010                   **Dimensions**  
 2011

---

2012                   Hadamard product of two matrices is the multiplication of the corresponding pair  
 2013                   of elements of the two matrices of the same dimension:

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

2014                   The following is a computer code for an example of Hadamard product.

```
2015 #R Hadamard product of two matrices
2016 #install.packages('matrixcalc')
2017 library(matrixcalc)
2018 A = matrix( c( 1, 2, 3, 4 ), nrow=2, byrow=TRUE )
2019 B = matrix( c( 2, 4, 6, 8 ), nrow=2, byrow=TRUE )
2020 hadamard(A, B)
2021 #      [,1] [,2]
2022 #[1,]    2     8
2023 #[2,]   18    32
```

2024                   An application of the Hadamard product is to calculate the value of each item of  
 2025                   a store, when the numbers of items for each good are stored as a matrix, and the  
 2026                   corresponding unit prices as another. The Hadamard product of the two matrices  
 2027                   yields a matrix of values for each kind of good.

2028                   **2.5.2 Jordan Product of Two Matrices of the Same Dimensions**  
 2029

---

2030                   Jordan product of two square matrices of the same dimensions is defined as follows:

2031

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

2032                   A numerical example is as follows.

```
2033 #R Jordan product of A and B
2034 A = matrix( c( 1, 2, 3, 4 ), nrow=2, byrow=TRUE )
2035 B = matrix( c( 2, 1, 2, 1 ), nrow=2, byrow=TRUE )
2036 (A %*% B + B %*% A)/2
2037 #      [,1] [,2]
2038 #[1,]    5.5  5.5
2039 #[2,]   9.5  7.5
```

2040                   Since matrix multiplication is not commutative,  $AB \neq BA$ , Jordan product is  
 2041                   the mean of  $AB$  and  $BA$ .

2042                   **2.5.3 Commutator of Two Matrices of the Same Dimensions**  
 2043

---

2044                   The commutator of two matrices of the same dimensions is defined as follows:

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

2045                   A numerical example is below.

```

2046 #R commutator of A and B
2047 #R commutator of A and B
2048 A = matrix( c( 1, 2, 3, 4 ), nrow=2, byrow=TRUE )
2049 B = matrix( c( 2, 1, 2, 1 ), nrow=2, byrow=TRUE )
2050 A%*%B - B%*%A
2051 # [,1] [,2]
2052 #[1,] 1 -5
2053 #[2,] 9 -1
2054 #install.packages('psych')
2055 library(psych)
2056 tr(A%*%B - B%*%A) #tr for trace
2057 #[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:

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

**2062 Theorem 2.1** *Shoda theorem: If  $\text{tr}(C)$  is zero, then  $C$  is a commutator matrix.*

2063 The proof of this theorem can made through constructing matrices  $A$  and  $B$  such  
2064 that

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

2065 Given  $C$ , it is similar to another trace zero matrix  $D$ , i.e., there exists an orthog-  
2066 onal matrix  $S$  such as

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

2067 Then, it turns out that

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

2068 where

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

2069 with

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

2070 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

2074 The outer product of two vectors  $a_m$  and  $b_{p \times 1}$  is defined as

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

2075 and is an  $m \times p$  matrix.

```

2076 #Outer product of two vectors
2077 a = 1:2
2078 b = 1:4
2079 a%o%b #outer product a_2-by-1 times t(b)_1-by-4
2080 # [1,] 1 2 3 4
2081 # [2,] 2 4 6 8
2083
2084 #Outer product of A_mn and B_nq
2085 A = matrix(1:4, ncol = 2)
2086 B = matrix(1:6, ncol = 3)
2087 A%o%B
2088 dim(A%o%B)
2089 #[1] 2 2 2 3
2090
2091 #Outer product of A_mn and B_pq
2092 A = matrix(1:4, ncol = 2)
2093 B = matrix(1:9, ncol = 3)
2094 A%o%B
2095 dim(A%o%B)
2096 #[1] 2 2 2 3

```

2097 This concept can be extended to the outer product of two matrices  $A$  and  $B$ , as  
2098 shown in this numerical example.

#### 2099 2.5.4.2 Cross Product of Two Vectors in 3D

2100 The cross product of two vectors  $x = (x_1, x_2, x_3)$  and  $y = (y_1, y_2, y_3)$  in 3-dimensional  
2101 space is also a 3-dimensional vector and defined as follows

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

2102 A numerical example is below.

```

2103 #R cross product and dot product
2104 library(pracma)
2105 x = 1:3
2106 y = 4:6
2107 cross(x, y)
2108 #[1] -3 6 -3
2109 dot(x, y)
2110 #[1] 32

```

2111 This example also shows a dot product of two vectors. The dot product can  
2112 be extended to n-dimensional space, but the cross product is limited to the 3-  
2113 dimensional space. The cross product is used often in physics and engineering,  
2114 particularly in mechanics and electromagnetism. The cross product of the force  
2115 vector with an arm vector is a torque vector, which points to the direction of a  
2116 screw when you drive it into a wall.

2117 The dot product in 3D space also has many physics and engineering applications,  
2118 such as the force's dot product with displacement vector is work.

2119      **2.5.4.3 Another Matrix Product Based on Outer Product of  
2120      Vectors**

2121      There is another type of product of two matrices using the outer product of column  
2122      and row vectors. We regard  $A_{mn} = [a[, i]], i = 1, 2, \dots, n$  consisting of  $n$  column  
2123      vectors, and  $B_{n \times q} = [b[j,], j = 1, 2, \dots, n]$  consisting of  $n$  row vectors. The outer  
2124      product is an  $mq$  matrix. defined as the

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

2125      **2.5.5 Kronecker Product of Two Matrices of the Same  
2126      Dimensions**  
2127

---

2128      Kronecker product of two matrices  $A_{mn}$  and  $B_{p \times q}$  is defined as the

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

2129      Two numerical examples are as follows:

```
2130 #R Kronecker product of two matrices
2131 #Kronecker product
2132 library(fastmatrix)
2133 A <- diag(1:2)
2134 B <- matrix(1:4, ncol = 2)
2135 kronecker.prod(A, B)
2136 # [,1] [,2] [,3] [,4]
2137 #[1,] 1 3 0 0
2138 #[2,] 2 4 0 0
2139 #[3,] 0 0 2 6
2140 #[4,] 0 0 4 8
2141
2142 # an example with vectors
2143 ones <- rep(1, 2)
2144 y <- 1:4
2145 kronecker.prod(ones, t(y)) # 2-by-4 matrix
2146 # [,1] [,2] [,3] [,4]
2147 #[1,] 1 2 3 4
2148 #[2,] 1 2 3 4
```

2149      **Theorem 2.2**  $(A \diamond B)^{-1} = A^{-1} \diamond B^{-1}$ . **Proof:**

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

2150      Here we have used the following equality

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

2151      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  $\oplus$ :

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

```

2154 #R direct sum of two matrices
2155 A = matrix(1:4, ncol = 2)
2156 B = matrix(1:6, ncol = 3)
2157 A1 = rbind(A, matrix(rep(0, 4), ncol = 2))
2158 B1 = rbind(matrix(rep(0, 6), ncol = 3), B)
2159 C = cbind(A1, B1) #= direct sum of A and B
2160 C
2161 #      [,1] [,2] [,3] [,4] [,5]
2162 #[1,]    1    3    0    0    0
2163 #[2,]    2    4    0    0    0
2164 #[3,]    0    0    1    3    5
2165 #[4,]    0    0    2    4    6
2166
2167 #Express the direct sum by Kronecker products
2168 kronecker.prod(diag(1,0), A) + kronecker.prod(diag(0,1), B)
2169 #      [,1] [,2] [,3] [,4] [,5]
2170 #[1,]    1    3    0    0    0
2171 #[2,]    2    4    0    0    0
2172 #[3,]    0    0    1    3    5
2173 #[4,]    0    0    2    4    6
2174
2175
2176

```

**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:

- 2188 (i) What are the typical differences between your clients?  
 2189 (ii) What are the temporal variation patterns that show the beginning of the  
 2190 effectiveness of your treatment?  
 2191 The eigenvalues and eigenvectors of  $C$  may help answer Question (i). The SVD  
 2192 in the next section may help answer Question (ii).  
 2193 We use the following numerical example to illustrate the procedure of visualiza-  
 2194 tion and interpretation of eigenvalues and eigenvectors for  $C$ .

2195 **2.8 Singular Value Decomposition**

---

2196 Previous section shows an eigenvector-eigenvalue decomposition of a symmetric  
 2197 square matrix. A similar decomposition can be made for a rectangular matrix. The  
 2198 decomposition using unit eigenvectors and eigenvalues for a general rectangular  
 2199 matrix is called the *singular value decomposition* (SVD). Singular value is another  
 2200 name of eigenvalue. Although the basic mathematical theory of SVD was devel-  
 2201 oped almost 200 years ago (Stewart 1993), the modern algorithm of efficient SVD  
 2202 computing was only developed by Gene H. Golub (1932-2007) and his colleagues in  
 2203 the 1970s. Now, SVD has become an important data analysis tool for every field.  
 2204 Climate science is not an exception. A space-time climate data matrix is often a  
 2205 rectangular matrix, since the number of sites is not likely to be equal to the number  
 2206 of temporal observations at those sites. We may denote a space-time climate data  
 2207 matrix by  $\mathbf{X}_{n \times m}$ , where  $n$  denotes the number of sites, and  $m$  is the total number  
 2208 of temporal observations. For  $\mathbf{X}_{n \times m}$ , we can also interpret  $n$  as the number of grid  
 2209 boxes of a climate model output, and  $m$  as the number of time steps.

2211 **2.8.1 SVD formula and a simple SVD example**

---

2212 If  $m \leq n$ , then matrix  $\mathbf{A}_{n \times m}$  has the following SVD decomposition

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

2213 Here,  $\mathbf{U}$  may be interpreted as a spatial matrix, consisting of  $m$  spatial orthonormal  
 2214 column vectors which are unit eigenvectors of  $\mathbf{A}\mathbf{A}^t$ ;  $\mathbf{V}$  may be interpreted as a  
 2215 temporal matrix, consisting of  $m$  temporal orthonormal column vectors which are  
 2216 unit eigenvectors of  $\mathbf{A}^t\mathbf{A}$ ; and  $\mathbf{D}$  is a diagonal matrix whose elements are the square  
 2217 root of the eigenvalues of  $\mathbf{A}\mathbf{A}^t$  and may be interpreted as standard deviations.  
 2218 Figure 2.3 may help you understand this formula. The diagonal elements of  $\mathbf{D}$   
 2219 are called singular values, and the column vectors of  $\mathbf{U}$  and  $\mathbf{V}$  are called singular  
 2220 vectors. Here, the word “singular” may be understood as special, or opposite to  
 2221 “general”, or distinguished, or being out of ordinary.

2222 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)

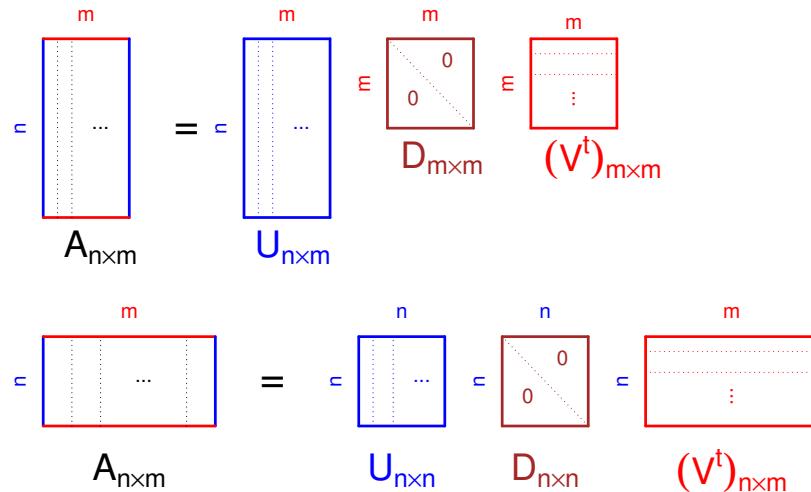


Fig. 2.3 Schematic diagrams of singular value decomposition (SVD):  $A = UDV^t$ .

2224 The following R code shows a simple SVD example of a  $2 \times 3$  matrix.

```

2225 #SVD example for a 2-by-3 matrix: R code
2226 A=matrix(c(-1,1,0,2,-2,3),nrow=2)
2227 A #Show the 2-by-3 matrix
2228 # [,1] [,2] [,3]
2229 #[1,] -1 0 -2
2230 #[2,] 1 2 3
2231 svdA=svd(A) #Compute the SVD of A and put the results in svdA
2232 svdA #Show SVD results: d, U, and V
2233 round(svdA$d, digits=2) #Show only the singular values
2234 #[1] 4.22 1.09
2235 round(svdA$u, digits=2) #Show only matrix U
2236 # [,1] [,2]
2237 #[1,] -0.48 0.88
2238 #[2,] 0.88 0.48
2239 round(svdA$v, digits=2) #Show only matrix V
2240 # [,1] [,2]
2241 #[1,] 0.32 -0.37
2242 #[2,] 0.42 0.88
2243 #[3,] 0.85 -0.29
2244 sqrt(eigen(A%*%t(A))$values)
2245 #[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])

```

2246

2247 The above computer code shows the following SVD results expressed in mathematical formulas below:

2248 (a) The vector form:

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

2250

2251 and

2252 (b) The matrix form:

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

2253 If we use only the first singular vectors to approximate  $A$  from the two triplets  
2254 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

```

2261 It is quite close to  $A$ .

```

2262
2263 #Data reconstruction by singular vectors: Python code
2264 np.round(DsvdA[0]*np.outer(UsvdA[:,0], VsvdA[:,0]),1)
2265 #array([[-0.7, -0.8, -1.7],
2266 #        [ 1.2,  1.5,  3.2]])

```

2262

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

```

2265 round(svdA$d[1]*svdA$u[,1]%^%t(svdA$v[,1]) +
2266   svdA$d[2]*svdA$u[,2]%^%t(svdA$v[,2]),
2267   digits =2)
2268 #[,1] [,2] [,3]
2269 #[1,] -1 0 -2
2270 #[2,] 1 2 3

```

2271

```

A1 = DsvdA[0]*np.outer(UsvdA[:,0], VsvdA[:,0])
A2 = DsvdA[1]*np.outer(UsvdA[:,1], VsvdA[:,1])
np.round(A1 + A2, 2)
#array([[-1.,  0., -2.],
#       [ 1.,  2.,  3.]])

```

2271

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

```

2272 #R plot schematic diagram of SVD
2273 setwd('/Users/sshen/LinAlg')
2274 setEPS() #Plot the figure and save the file
2275 postscript("fig0503.eps", width = 11)
2276 par(mar=c(0,0,0,0))
2277 plot(200, axes = FALSE,
2278       xlab = "", ylab = "",
2279       xlim = c(-3,28), ylim = c(-3,16))
2280 text(13,15.5, cex=2.2,
2281       bquote("SVD:" ~ A == UDV^t ~ "when n > m or n < m"))
2282 #Space-time data matrix A when n>m
2283 segments(x0 = c(0,0,3,3),
2284            y0 = c(6,12,12,6) +1,
2285            x1 = c(0,3,3,0),
2286            y1 = c(12,12,6,6) +1,
2287            col = c('blue','red','blue','red')), lwd = 3)
2288 segments(x0 = c(0.5,1.0),
2289            y0 = c(6,6)+1,
2290            x1 = c(0.5,1.0),
2291            y1 = c(12,12)+1,
2292            lwd = 1.3, lty = 3)
2293 text(-.8, 9+1, 'n', srt=90, col ='blue', cex = 1.4)
2294 text(1.5, 12.8+1, 'm', col = 'red', cex = 1.4)
2295 text(2.0, 9+1, '...', cex = 1.4)
2296 text(2, 5+1, bquote(A[n %*% m]), cex = 2.5)
2297 text(5, 9+1, '=', cex = 3)
2298 #Spatial matrix U

```

```

2301 segments(x0 = c(7,7,10,10),
2302           y0 = c(6,12,12,6)+1,
2303           x1 = c(7,10,10,7),
2304           y1 = c(12,12,6,6)+1,
2305           col = c('blue','blue','blue','blue'), lwd =3)
2306 segments(x0 = c(7.5,8),
2307           y0 = c(6,6)+1,
2308           x1 = c(7.5,8),
2309           y1 = c(12,12)+1,
2310           lwd =1.3, lty = 3, col = 'blue')
2311 text(6.2, 9+1, 'n', srt=90, col = 'blue', cex = 1.4)
2312 text(8.5, 12.8+1, 'm', col = 'red', cex = 1.4)
2313 text(9, 9+1, '...', cex = 1.4, col='blue')
2314 text(8.7, 5.0+1, bquote(U[n%*%m]), cex = 2.5, col= 'blue')
2315 #Singular value diagonal matrix D
2316 segments(x0 = c(12,12,15,15),
2317           y0 = c(9,12,12,9)+1,
2318           x1 = c(12,15,15,12),
2319           y1 = c(12,12,9,9)+1,
2320           col = c('brown','brown','brown','brown'), lwd =3)
2321 segments(x0 = 12, y0 = 12+1, x1 = 15, y1 = 9+1, lty=3,
2322           col = c('brown'), lwd =1.3)#diagonal line
2323 text(11.2, 10.5+1, 'm', srt=90, col ='red', cex = 1.4)
2324 text(13.5, 12.8+1, 'm', col = 'red', cex = 1.4)
2325 text(14.1, 11.3+1, '0', col = 'brown', cex = 1.4)
2326 text(12.9, 10.0+1, '0', col = 'brown', cex = 1.4)
2327 text(13.9, 8.0+1, bquote(D[m%*%m]), cex = 2.5, col='brown')
2328 #Temporal matrix V
2329 segments(x0 = c(17,17,20,20),
2330           y0 = c(9,12,12,9)+1,
2331           x1 = c(17,20,20,17),
2332           y1 = c(12,12,9,9)+1,
2333           col = c('red','red','red','red'), lwd =3)
2334 segments(x0 = c(17,17),
2335           y0 = c(11.5,10.8)+1,
2336           x1 = c(20,20),
2337           y1 = c(11.5,10.8)+1,
2338           col = c('red','red'), lty=3, lwd =1.3)
2339 text(16.2, 10.5+1, 'm', srt=90, col ='red', cex = 1.4)
2340 text(18.5, 12.5+1, 'm', col = 'red', cex = 1.4)
2341 text(19.5, 8+1, bquote((V^t)[m%*%m]), cex = 2.5, col='red')
2342 text(18.5, 10+1, '...', col='red', srt=90, cex =1.4)
2343 # Space-time data matrix B when n < m
2344 segments(x0 = c(0,0,6,6),
2345           y0 = c(0,3,3,0),
2346           x1 = c(0,6,6,0),
2347           y1 = c(3,3,0,0),
2348           col = c('blue','red','blue','red'), lwd =3)
2349 segments(x0 = c(1,2,5),
2350           y0 = c(0,0,0),
2351           x1 = c(1,2,5),
2352           y1 = c(3,3,3),
2353           lwd =1.3, lty = 3)
2354 text(-0.8, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2355 text(3, 3.8, 'm', col = 'red', cex = 1.4)
2356 text(3.5, 1.5, '...', cex = 1.4)

```

```

2357 text(3, -1.5, bquote(A[n%*%m]), cex = 2.5)
2358 text(8, 1.5, '=', cex = 3)
2359 #Spatial matrix U
2360 segments(x0 = c(11,11,14,14),
2361           y0 = c(0,3,3,0),
2362           x1 = c(11,14,14,11),
2363           y1 = c(3,3,0,0),
2364           col = c('blue','blue','blue','blue'), lwd =3)
2365 segments(x0 = c(11.5,12.2),
2366           y0 = c(0,0),
2367           x1 = c(11.5,12.2),
2368           y1 = c(3,3),
2369           lwd =1.3, lty = 3, col = 'blue')
2370 text(10.2, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2371 text(12.5, 3.8, 'n', col = 'blue', cex = 1.4)
2372 text(13.2, 1.5, '...', cex = 1.4, col='blue')
2373 text(12.5, -1.5, bquote(U[n%*%n]), cex = 2.5, col= 'blue')
2374 #Singular value diagonal matrix D
2375 segments(x0 = c(16,16,19,19),
2376           y0 = c(0,3,3,0),
2377           x1 = c(16,19,19,16),
2378           y1 = c(3,3,0,0),
2379           col = c('brown','brown','brown','brown'), lwd =3)
2380 segments(x0 = 16, y0 = 3, x1 = 19, y1 = 0, lty=3,
2381           col = c('brown')), lwd =1.3)#diagonal line
2382 text(15.2, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2383 text(17.5, 3.8, 'n', col = 'blue', cex = 1.4)
2384 text(18.1, 2.3, '0', col = 'brown', cex = 1.4)
2385 text(16.9, 1.0, '0', col = 'brown', cex = 1.4)
2386 text(17.5, -1.5, bquote(D[n%*%n]), cex = 2.5, col='brown')
2387 #Temporal matrix V
2388 segments(x0 = c(21,21,27,27),
2389           y0 = c(0,3,3,0),
2390           x1 = c(21,27,27,21),
2391           y1 = c(3,3,0,0),
2392           col = c('red','red','red','red'),
2393           lwd =3)
2394 segments(x0 = c(21,21),
2395           y0 = c(2.5,1.8),
2396           x1 = c(27,27),
2397           y1 = c(2.5,1.8),
2398           col = c('red','red'), lty=3, lwd =1.3)
2399 text(20.2, 1.5, 'n', srt=90, col ='blue', cex = 1.4)
2400 text(24, 3.8, 'm', col = 'red', cex = 1.4)
2401 text(24, -1.5, bquote((V^t)[n%*%m]), cex = 2.5, col='red')
2402 text(24, 1, '...', col='red', srt=90, cex =1.4)
2403 dev.off()

```

```

#Python plot Fig. 5.3: Schematic diagram of SVD
import matplotlib.patches as patches
import numpy as np
import pylab as pl
from matplotlib import collections as mc
lines = [(0, 7), (0, 13)], [(0, 13), (3, 13)],
          [(3, 13), (3, 7)], [(3, 7), (0, 7)]]
c = np.array(['b', 'r', 'b', 'r'])
lc = mc.LineCollection(lines, colors=c, linewidths=3)
fig, ax = pl.subplots()
ax.set_xlim([-3, 28])
ax.set_ylim([-3, 16])
ax.add_collection(lc)
ax.margins(0.1)
plt.plot([0.5, 0.5], [7, 13],
          linestyle='dotted', color = 'k')
plt.plot([1, 1], [7, 13],
          linestyle='dotted', color = 'k')
plt.text(13, 15.5,
         r'SVD: $A = UDV^T$ when $n > m$ or $n < m$',
         fontsize = 30)
plt.text(-1.2, 10, 'n', color = 'blue',
         fontsize = 25, rotation=90)
plt.text(1.1, 13.3, 'm', color = 'red',
         fontsize = 25, rotation=0)
plt.text(0.0, 5.5, r'$A_{\{n \times m\}}$',
         fontsize = 35, rotation=0)
plt.text(1.5, 10, '...', color = 'black',
         fontsize = 25, rotation=0)
plt.axis('off')
plt.show()

```

2404

2405 This Python code generates the top-left rectangular box in Fig. 5.3. The re-  
 2406 maining code for other boxes is highly repetitive and can be found from the book  
 2407 website.

2408

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

2409

2410 The Southern Oscillation Index (SOI) is an indicator for El Niño or La Niña. It is  
 2411 computed as the difference of sea level pressure (SLP) of Tahiti ( $17.75^{\circ}\text{S}$ ,  $149.42^{\circ}\text{W}$ )  
 2412 minus that of Darwin ( $12.46^{\circ}\text{S}$ ,  $130.84^{\circ}\text{E}$ ). An SVD analysis of the SLP data can  
 2413 substantiate this calculation formula.

2414

2415 The following shows the data matrix of the standardized SLP anomalies of Tahiti  
 2416 and Darwin from 2009 to 2015, and its SVD:

```

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

```

```

2419 Pda<-read.table("data/PSTANDdarwin.txt", header=F)
2420 dim(Pda)
2421 #[1] 65 13 #Monthly Darwin data from 1951-2015
2422 pdaDec<-Pda[,13] #Darwin Dec standardized SLP anomalies data
2423 Pta<-read.table("data/PSTANDtahiti.txt", header=F)
2424 ptaDec=Pta[,13] #Tahiti Dec standardized SLP anomalies
2425 ptada1 = cbind(pdaDec, ptaDec) #space-time data matrix
2426
2427 #Space-time data format
2428 ptada = t(ptada1[59:65,]) #2009-2015 data
2429 colnames(ptada)<-2009:2015
2430 rownames(ptada)<-c("Darwin", "Tahiti")
2431 ptada #6 year of data for two stations
2432 # 2009 2010 2011 2012 2013 2014 2015
2433 #Darwin 0.5 -2.3 -2.2 0.3 0.3 0.1 -0.4
2434 #Tahiti -0.7 2.5 1.9 -0.7 0.4 -0.8 -1.3
2435 svdptd = svd(ptada) #SVD for the 2-by-6 matrix
2436 U=round(svdptd$u, digits=2)
2437 U
2438 #[1,] -0.66 0.75
2439 #[2,] 0.75 0.66
2440 D=round(diag(svdptd$d), digits=2)
2441 D
2442 #[1,] 4.7 0.00
2443 #[2,] 0.0 1.42
2444 V =round(svdptd$v, digits=2)
2445 t(V)
2446 #[1,] -0.18 0.72 0.61 -0.15 0.02 -0.14 -0.15
2447 #[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] ...]

```

2449 One can verify that

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

2450 approximately recovers the original data matrix.

2451 The first column vector  $(-0.66, 0.75)$  of the spatial pattern matrix  $\mathbf{U}$  may be  
2452 interpreted to be associated with the Southern Oscillation Index (SOI), which puts  
2453 a negative weight  $-0.66$  on Darwin, and a positive weight  $0.75$  on Tahiti. The  
2454 weighted sum is approximately equal to the difference of Tahiti's SLP minus that  
2455 of Darwin, which the definition of SOI. The index measures large scale ENSO  
2456 dynamics of the tropical Pacific (Trenberth 2020). The magnitude of the vector  
2457  $(-0.66, 0.75)$  is approximately one, because  $\mathbf{U}$  is a unitary matrix. The correspond-  
2458 ing singular vector has a distinctly large value  $0.72$  in December 2010, which was  
2459 a strong La Niña month. In this month, the Darwin has a strong negative SLP  
2460 anomaly, while Tahiti has a strong positive SLP anomaly. This situation enhanced  
2461 the easterly trade winds in the tropical Pacific and caused abnormally high precip-  
2462 itation in Australia in the 2010-2011 La Niña period.

2463 The second column vector  $(0.75, 0.66)$  of the spatial pattern matrix  $\mathbf{U}$  also has  
2464 climate implications. The weighted sum with two positive equal weights measures  
2465 the tropical Pacific dynamics of small scales (Trenberth 2020).

## 2466 2.10 Chapter summary

---

2467 A matrix can be regarded as a 2-dimensional  $n \times m$  rectangular array of numbers  
2468 or symbols, or as  $m$  column vectors, or as  $n$  row vectors. This may be interpreted as  
2469 climate data at  $n$  locations with  $m$  time steps. Many mathematical properties of a  
2470 matrix of climate data have climate interpretations. For example, SVD decomposes  
2471 a space-time climate data matrix into an orthogonal spatial pattern matrix, an  
2472 orthogonal temporal pattern matrix, and a diagonal “energy-level” matrix that  
2473 measures the standard deviation of the temporal pattern:

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

2475 The column vectors of the spatial matrix  $\mathbf{U}$  are spatial singular vectors, also called  
2476 EOFs, while those of the temporal matrix  $\mathbf{V}$  are temporal singular vectors, also  
2477 called PCs. The first a few EOFs often have climate dynamic interpretations, such  
2478 as El Niño Southern Oscillation (ENSO). If EOF1 corresponds to El Niño and  
2479 shows some typical ENSO properties, such as the opposite signs of SLP anomalies  
2480 of Darwin and Tahiti, then PC1 shows a temporal pattern, e.g., the extreme values  
2481 of PC1 indicating both the occurrence time and the strength of El Niño. The  
2482 diagonal elements of matrix  $\mathbf{D}$  are singular values, also known as eigenvalues.

2483 An eigenvector  $\mathbf{v}$  of a square matrix  $\mathbf{C}$  is such a special vector that  $\mathbf{C}$ 's action  
2484 on  $\mathbf{v}$  does not change its orientation, i.e.,  $\mathbf{C}\mathbf{v}$  is parallel to  $\mathbf{v}$ . This statement

2485 implies the existence of a scalar  $\lambda$ , called eigenvalue (also known as singular value  
2486 or characteristic value), such that

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

2487 We have also discussed the matrix method of solving a system of linear equations  
2488

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

2489 linear independence of vectors, linear transform, and other basic matrix methods.  
2490 These methods are useful for the chapters on covariance, EOFs, spectral analysis,  
2491 regression analysis, and machine learning. You may focus on the computing meth-  
2492 ods and computer code of the relevant methods. The mathematical proofs of this  
2493 chapter, although helpful for exploring new mathematical methods, are not neces-  
2494 sarily needed to read the other chapters of this book. If you are interested in an  
2495 in-depth mathematical exploration of matrix theory, you may wish read the books  
2496 by Horn and Johnson (1985) and Strang (2016).

## References and Further Readings

- 2498 [1] Golub, G.H. and C. Reinsch, 1970: Singular value decomposition and least  
 2499 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.

- 2500  
 2501 [2] Horn, R. A., and C. R. Johnson, 1985: *Matrix Analysis*, Cambridge University  
 2502 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.

- 2503  
 2504 [3] Strang, G., 2016: *Introduction to Linear Algebra*. 5th edition, Wellesley-  
 2505 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.

- 2506  
 2507 [4] Stewart, G.W., 1993: On the early history of the singular value decomposition.  
 2508 *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.

- 2509  
 2510 [5] Trenberth, K., and National Center for Atmospheric Research Staff (Eds), 2020:  
 2511 The Climate Data Guide: Southern Oscillation Indices: Signal, Noise and Tahiti-  
 2512 Darwin SLP (SOI). Retrieved from

2513 <https://climatedataguide.ucar.edu/climate-data/southern-oscillation-indices-signal-noise-and-tahitidarwin-slp-soi>

2515 This site describes the optimal indices for large and small scale dynamics.

- 2516 [6] Tucker, A., 1993: The growing importance of linear algebra in undergraduate  
2517 mathematics. *College Mathematics Journal* 24, 3-9.

2518 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.

## 2519 Exercises

- 2520
- 2521 **2.1** Write a computer code to  
 2522 (a) Read the NOAAGlobalTemp data file, and  
 2523 (b) Generate a  $4 \times 8$  space-time data matrix for the December mean surface  
 2524 air temperature anomaly data of four grid boxes and eight years.  
 2525 *Hint: You may find the NOAA Global Surface Temperature (NOAAGlobal-  
 2526 Temp) dataset online. You can use either netCDF format or CSV format.*

- 2527 **2.2** Write a computer code to find the inverse of the following matrix

```
2528 # [,1] [,2] [,3]
2529 #[1,] 1.7 -0.7 1.3
2530 #[2,] -1.6 -1.4 0.4
2531 #[3,] -1.5 -0.3 0.6
```

- 2532 **2.3** Write a computer code to solve the following linear system of equations

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

2533 where

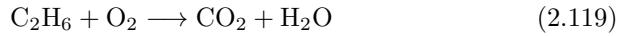
$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 0 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \quad (2.117)$$

- 2534 **2.4** The following equation

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

2535 has infinitely many solutions, and cannot be directly solved by a simple com-  
 2536 puter command, such as `solve(A, b)`.  
 2537 (a) Show that the three row vectors of the coefficient matrix are not linearly  
 2538 independent.  
 2539 (b) Because of the dependence, the linear system has only two independent  
 2540 equations. Thus, reduce the linear system into two equations by treating  $x_3$   
 2541 as an arbitrary value while treating  $x_1$  and  $x_2$  as variables.  
 2542 (c) Solve the two equations for  $x_1$  and  $x_2$  and express them in terms of  $x_3$ .  
 2543 The infinite possibilities of  $x_3$  imply infinitely many solutions of the original  
 2544 system.

2545 **2.5** Ethane is a gas similar to the greenhouse gas methane and can burn with  
 2546 oxygen to form carbon dioxide and water:



2547 Given two ethane molecules, how many molecules of oxygen, carbon dioxide  
 2548 and water should be in order for this chemical reaction equation to be bal-  
 2549 anced?

2550 *Hint: Assume  $x, y$  and  $z$  molecules of oxygen, carbon dioxide and water, and  
 2551 use the balance of the number of atoms of carbon, hydrogen and oxygen to  
 2552 form linear equations. Solve the system of linear equations.*

2553 **2.6** Carry out the same procedure as the previous problem but for the burning of  
 2554 methane  $\text{CH}_4$ .

2555 **2.7** (a) Use matrix multiplication to show that the vector

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

2556 is not an eigenvector of the following matrix

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

2557 noindent (b) Find all the unit eigenvectors of matrix  $\mathbf{A}$  in Part (a).

2558 **2.8** Use hand-calculation to compute the matrix multiplication of  $\mathbf{UDV}^t$  where  
 2559 the data of relevant matrices are given by the following R output:

```
2560 A=matrix(c(1,-1,1,1),nrow=2)
2561 A
2562 #[1,]    1    1
2563 #[2,]   -1    1
2564 svd(A)
2565 #\$d
2566 #[1] 1.414214 1.414214
2567 #\$u
2568 #[1,] -0.7071068 0.7071068
2569 #[2,]  0.7071068 0.7071068
2570 #\$v
2571 #[1,]   -1    0
2572 #[2,]    0    1
```

- 2573    **2.9** Use computer to find the matrices  $\mathbf{U}$ ,  $\mathbf{D}$  and  $\mathbf{V}$  of the SVD for the following  
2574    data matrix

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

- 2575    **2.10** Use the first singular vectors to approximately reconstruct the data matrix  
2576     $\mathbf{A}$  using

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

2577    Describe the goodness of the approximation using text, limited to 20 - 100  
2578    words.

- 2579    **2.11** For the following data matrix

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

- 2580    (a) Use computer to find the eigenvectors and eigenvalues of matrix  $\mathbf{A}\mathbf{A}^t$ .
- 2581    (b) Use computer to find the eigenvectors and eigenvalues of matrix  $\mathbf{A}^t\mathbf{A}$ .
- 2582    (c) Use computer to calculate SVD of  $\mathbf{A}$ .
- 2583    (d) Compare the singular vectors and singular values in Step (c) with the  
2584    eigenvalues and eigenvectors computed in Steps (a) and (b). Use text to discuss  
2585    your comparison.
- 2586    (e) Use computer to verify that the column vectors of  $\mathbf{U}$  and  $\mathbf{V}$  in the SVD  
2587    of Step (c) are orthonormal to each other.

- 2588    **2.12** Make an SVD analysis of the December standardized anomalies data of sea  
2589    level pressure (SLP) at Darwin and Tahiti from 1961 to 2010. You can find  
2590    the data from the Internet or from the website of this book.
- 2591    (a) Write a computer code to organize the data into a  $2 \times 50$  space-time data  
2592    matrix.
  - 2593    (b) Make the SVD calculation for this space-time matrix.
  - 2594    (c) Plot the first singular vector in  $\mathbf{V}$  against time 1961 to 2010 as a time  
2595    series curve, which is called the first principal component, denoted by PC1.
  - 2596    (d) Interpret the first singular vector in  $\mathbf{U}$ , which is called the first empirical  
2597    orthogonal function (EOF1), as weights of Darwin and Tahiti stations.
  - 2598    (e) Check the historical El Niño events between 1961 and 2010 from the  
2599    Internet, and interpret the extreme values of PC1.

- 2600    **2.13** Plot PC2 against time from the SVD analysis in the previous problem. Discuss  
2601    the singular values  $\lambda_1$  and  $\lambda_2$  and interpret PC2, in reference to the description  
2602    entitled “The Climate Data Guide: Southern Oscillation Indices: Signal, Noise  
2603    and Tahiti/Darwin SLP (SOI)” by Kevin Trenberth (2020).

- 2604    **2.14** Make an SVD analysis similar to the previous two problems for the January  
2605    standardized SLP anomalies data at Darwin and Tahiti from 1961 to 2010.

- 2606   **2.15** Make an SVD analysis similar to the previous problem for the monthly stan-  
 2607   dardized SLP anomalies data at Darwin and Tahiti from 1961 to 2010. This  
 2608   problem includes anomalies for every month. The space-time data is a  $2 \times 600$   
 2609   matrix.
- 2610   **2.16** For the observed data of the monthly surface air temperature at five stations  
 2611   of your interest from January 1961 to December 2010, form a space-time data  
 2612   matrix, compute the SVD of this matrix, and interpret your results from the  
 2613   perspective of climate science. Please plot PC1, PC2, and PC3 against time.  
 2614   You may find your data from the Internet, such as the NOAA Climate Data  
 2615   Online website <https://www.ncdc.noaa.gov/cdo-web>
- 2616   **2.17** Do the same analysis as the previous problem, but for the monthly precipita-  
 2617   tion data at the same stations in the same time period.
- 2618   **2.18** For an Reanalysis dataset, make an SVD analysis similar to the previous  
 2619   problem for the monthly surface temperature data over 10 grid boxes of your  
 2620   choice in the time period of 1961-2010. The space-time data is a  $10 \times 600$   
 2621   matrix. You can use your preferred Reanalysis dataset and download the data  
 2622   from the Internet, such as NCEP/NCAR Reanalysis, and ECMWF ERA.
- 2623   **2.19** Let  $\mathbf{C} = \mathbf{A}\mathbf{A}^t$  and  $\mathbf{A}$  is any real-valued rectangular matrix. Show that  
 2624   (a) the eigenvalues of  $\mathbf{C}$  are non-negative;  
 2625   (b) if  $\mathbf{v} = \mathbf{A}^t\mathbf{u}$ , then  $\mathbf{v}$  is a eigenvector of  $\mathbf{C}^t = \mathbf{A}^t\mathbf{A}$ .
- 2626   **2.20** Given that

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

2627   write down the second order polynomial corresponding to the following matrix  
 2628   expression

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

- 2629   This is known as the quadratic form of the matrix  $\mathbf{A}\mathbf{A}^t$ .  
 2630   **2.21** If  $\mathbf{x}$  is a unit vector, use calculus to find the maximum value of  $P(x_1, x_2)$  in  
 2631   the previous problem. How is your solution related to the eigenvalue of  $\mathbf{A}\mathbf{A}^t$ .

2632  
**3**

## Matrix Applications to Machine Learning

2634

2635 Machine learning (ML) is a branch of science that uses data and algorithms to  
2636 mimic how human beings learn. The accuracy of the ML results can be gradually  
2637 improved based on new training data and algorithm update. For example, a baby  
2638 learns how to pick an orange from a fruit plate containing apples, bananas and  
2639 oranges. Another baby learns how to sort out different kinds of fruits from a basket  
2640 into three categories without naming the fruits. Then, how does ML work? It is  
2641 basically a decision process for clustering, classification, or prediction, based on  
2642 the input data, decision criteria, and algorithms. It does not stop here. It further  
2643 validates the decision results and quantifies errors. The errors and the updated data  
2644 will help update the algorithms and improve the results.

2645 ML has recently become a very popular method in climate science due to the  
2646 availability of powerful and convenient resources of computing. It has been used to  
2647 predict weather and climate, and to develop climate models. This chapter is a brief  
2648 introduction of ML and provides basic ideas and examples. Our materials will help  
2649 readers understand and improve the more complex ML algorithms used in climate  
2650 science, so that they can go a step beyond only applying the ML software packages  
2651 as a black box. We also provide R and Python codes for some basic ML algorithms,  
2652 such as K-means for clustering, support vector machine for the maximum separation  
2653 of sets, random forest of decision trees for classification and regression, and neural  
2654 network training and predictions.

2655 Artificial intelligence (AI) allows computers to automatically learn from past  
2656 data without human programming, which enables a machine to learn and to have  
2657 intelligence. Machine learning is a subset of AI. Our chapter here focuses on ML,  
2658 not the general AI.

2659

### 3.1 K-means clustering

2660

2661 A few toddlers at a daycare center may learn how to grab a few candies near  
2662 themselves. It can be a point of fighting for a candy at a location that is not  
2663 obviously closer to one than another. K-means method can help divide the candies  
2664 among the toddlers in a fair way.

2665 Based on the historical weather data over a country, can ML decide the climate  
2666 regimes for the country? Can ML determine the ecoregions of a country? Can ML

2667 define the regimes of wild fire over a region? The K-means clustering method can  
 2668 be useful to answering these questions.

### 2669 2670 3.1.1 K-means setup and trivial examples

2671 The aim of the K-means clustering is to divide  $N$  points into  $K$  clusters so that  
 2672 the total within cluster sum of squares (tWCSS) is minimized. Here we use 2D  
 2673 data points to describe tWCSS and the K-means algorithm, although the K-means  
 2674 method can be formulated in higher dimensions. We regard the data  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$   
 2675 as the 2D coordinates of  $N$  points. For example, we treat  $\mathbf{x}_1 = (1.2, -7.3)$  as obser-  
 2676 vational data. Assume that these points can be divided into  $K$  clusters  $(C_1, C_2, \dots, C_K)$ ,  
 2677 where  $K$  is subjectively given by you, the user who wishes to divide the  $N$  points  
 2678 into  $K$  clusters by the K-means method. Let  $\mathbf{x} \in C_i$ , i.e., points within cluster  $C_i$ .  
 2679 The number of points in cluster  $C_i$  is unknown and is to be determined by the  
 2680 K-means algorithm. The total WCSS is defined by the following formula

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

2681 where

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

2682 is the mean of the points within the cluster  $C_i$ , as  $K_i$  is the number of points in  $C_i$ .  
 2683 Thus, we have K means, the name of the K-means method. That in the parentheses  
 2684 is the within cluster sum of squares (WCSS)

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

2685 This is defined for each cluster, and tWCSS is the sum of these  $\text{WCSS}_i$  for all the  
 2686  $K$  clusters.

2687 Some publications use WCSS or Tot WCSS, instead of tWCSS, to express the  
 2688 right hand side of Eq. (3.1). Be careful when reading literatures concerning the  
 2689 WCSS definition. Some computer software for K-means may even use a different  
 2690 tWCSS definition.

2691 Our K-means computing algorithm is to minimize the tWCSS by optimally or-  
 2692 ganizing the  $N$  points into  $K$  clusters. This algorithm can assign each data point  
 2693 to a cluster. If we regard  $\boldsymbol{\mu}_i$  as the centroid or center of cluster  $C_i$ , we may say that  
 2694 the points in cluster  $C_i$  may have some kind of similarity, such as similar climate or  
 2695 similar ecological characteristics, based on certain criteria. Below we use two simple  
 2696 cases ( $N = 2$  and  $N = 3$ ) to illustrate the K-means method and its solutions.

2697 (i) Case  $N = 2$ :

2698

2699 When  $N = 2$ , if we assume  $K = 2$ , then

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

2700 and

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

2701 This is the global minimum for  $K = 2$ , and

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

2702 is the unique solution. Of course, this is a trivial solution and bears no meaning.

2703 When  $N = 2$ , if we assume  $K = 1$ , then

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

2704 and

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

2705 This is the unique minimum of tWCSS, and the solution is also trivial.

2706 **(ii) Case  $N = 3$ :**

2707 When  $N = 3$ , if we assume  $K = 2$ , then there are three different cases of clustering

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

2710 Here, the case

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

2711 and

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

2712 are considered the same, since their difference is only a matter which cluster is  
2713 called  $C_1$  or  $C_2$ .

2714 For case  $C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}$ ,  $C_2 = \{\mathbf{x}_3\}$ , we have

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

2715 and

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

2716 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  
2717 be computed in a similar way. Then, the smallest tWCSS from the three determines  
2718 the final solution, which is one of the three cases of clustering. A numerical example  
2719 is given below.

2720

---

2721 **Example 3.1** Find the two K-means clusters for the following three points

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

2722 This is a K-means problem with  $N = 3$  and  $K = 2$ . The K-means clusters can be  
 2723 found by a computer command, e.g., `kmeans(mydata, 2)` in R. Figure 3.1 shows  
 2724 the two K-means clusters with their centers at

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

2725 and

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

2726 Points  $P_1$  and  $P_2$ , indicated by red dots, are assigned to cluster  $C_1$  whose center is  
 2727 marked by a red cross. Point  $P_3$ , indicated by the sky blue dot, is assigned to cluster  
 2728  $C_2$ , whose center is marked with a blue cross and obviously overlaps with  $P_3$ . This  
 2729 result agrees with our intuition since  $P_1$  and  $P_2$  are together. The three points can  
 2730 have two other possibilities of combination  $C_1 = \{P_1, P_3\}$  and  $C_1 = \{P_2, P_3\}$ . The  
 2731 case  $C_1 = \{P_1, P_3\}$  leads to  $\text{tWCSS} = 5.125$ , and  $C_1 = \{P_2, P_3\}$  to  $\text{tWCSS} = 3.625$ .  
 2732 Both tWCSS are greater than 0.5. Thus, these two combinations should not be a  
 2733 solution of the K-means clustering. These numerical results may be produced by  
 2734 the following computer codes.

```
2735 #tWCSS calculation for N = 3 and K = 2
2736 mydata <- matrix(c(1, 1, 2, 1, 3, 3.5),
2737   nrow = N, byrow = TRUE)
2738 x1 = mydata[1, ]
2739 x2 = mydata[2, ]
2740 x3 = mydata[3, ]
2741
2742 #Case C1 = (P1, P2)
2743 c1 = (mydata[1, ] + mydata[2, ])/2
2744 c2 = mydata[3, ]
2745 tWCSS = norm(x1 - c1, type = '2')^2 +
2746   norm(x2 - c1, type = '2')^2 +
2747   norm(x3 - c2, type = '2')^2
2748 tWCSS
2749 #[1] 0.5
2750
2751 #Case C1 = (P1, P3)
2752 c1 = (mydata[1, ] + mydata[3, ])/2
2753 c2 = mydata[2, ]
2754 norm(x1 - c1, type = '2')^2 +
2755   norm(x3 - c1, type = '2')^2 +
2756   norm(x2 - c2, type = '2')^2
2757 #[1] 5.125
2758
2759 #Case C1 = (P2, P3)
2760 c1 = (mydata[2, ] + mydata[3, ])/2
2761 c2 = mydata[1, ]
2762 norm(x2 - c1, type = '2')^2 +
2763   norm(x3 - c1, type = '2')^2 +
```

```
2764     norm(x1 - c2, type = '2')^2
2765 #[1] 3.625
2766
2767 #The case C1 = (P1, P2) can be quickly found by
2768 kmeans(mydata, 2)
2769 #Clustering vector:
2770 #[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

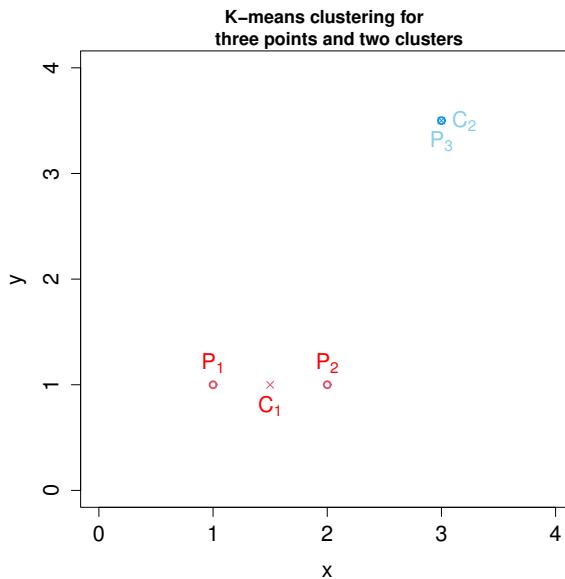
# case C1 = (P1, P2) can be quickly found by
kmeans = KMeans(n_clusters=2).fit(mydata)

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 ( $K = 2$ ).

2772 Of course, this problem is trivial and can even be solved analytically by hand.  
 2773 The computer solution of the K-means clustering is necessary when more points  
 2774 are involved and when it is not obvious for a point to belong to a specific cluster.  
 2775

2776 Figure 3.1 may be generated by the following computer code.

```

2777 # R plot Fig. 9.1: K-means for N = 3 and K = 2
2778 setwd("~/climstats")
2779 N = 3 #The number of data points
2780 K = 2 #Assume K clusters
2781 mydata = matrix(c(1, 1, 2, 1, 3, 3.5),
2782                   nrow = N, byrow = TRUE)
2783 Kclusters = kmeans(mydata, K)
2784 Kclusters #gives the K-means results,
2785 #e.g., cluster centers and WCSS
2786 #Cluster means:
2787 #[,1] [,2]
2788 #1 1.5 1.0
2789 #2 3.0 3.5
2790 #Within cluster sum of squares by cluster:
2791 # [1] 0.5 0.0
2792 Kclusters$centers
2793 par(mar = c(4,4,2.5,0.5))
2794 plot(mydata[,1], mydata[,2], lwd = 2,
2795       xlim = c(0, 4), ylim = c(0, 4),
2796       xlab = 'x', ylab = 'y', col = c(2, 2, 4),
2797       main = 'K-means clustering for
2798       three points and two clusters',

```

```

2799      cex.lab = 1.4, cex.axis = 1.4)
2800 points(Kclusters$centers[,1], Kclusters$centers[,2],
2801   col = c(2, 4), pch = 4)
2802 text(1.5, 0.8, bquote(C[1])), col = 'red', cex = 1.4)
2803 text(3.2, 3.5, bquote(C[2])), col = 'skyblue', cex = 1.4)
2804 text(1, 1.2, bquote(P[1])), cex = 1.4, col = 'red')
2805 text(2, 1.2, bquote(P[2])), cex = 1.4, col = 'red')
2806 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\u2022three\u2022points\u2022and\u2022two\u2022clusters", 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()

```

2807

### 3.1.2 A K-means algorithm

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

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

$$d_i = |\mathbf{x} - \mu_i^{(1)}| \quad i = 1, 2, \dots, K. \quad (3.17)$$

2815 Assign the point  $\mathbf{x}$  to the cluster with the smallest distance. Then, assign the  
 2816 next point to a cluster until every point is assigned. Thus, the initial  $K$  clusters  
 2817 are formed:  $C_i^{(1)}$  and each  $C_i^{(1)}$  containing  $K_i^{(1)}$  data points ( $i = 1, 2, \dots, K$ ).  
 2818 • Update Step: Compute the centroids of the initial  $K$  clusters

$$\boldsymbol{\mu}_i^{(2)} = \frac{1}{K_i^{(1)}} \sum_{\mathbf{x} \in C_i^{(1)}} \mathbf{x} \quad (3.18)$$

2819 and compute tWCSS<sup>(2)</sup>

$$\text{tWCSS}^{(2)} = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i^{(2)}} |\mathbf{x} - \boldsymbol{\mu}_i^{(2)}|^2 \quad (3.19)$$

2820 These  $\boldsymbol{\mu}_i^{(2)}$  ( $i = 1, 2, \dots, K$ ) are regarded as the updated centroids. We then  
 2821 repeat the Assignment Step to assign each point to a cluster again according  
 2822 to the rule of the nearest distance. Next, we compute  $\boldsymbol{\mu}_i^{(3)}$  ( $i = 1, 2, \dots, K$ ),  
 2823 and tWCSS<sup>(3)</sup>

2824 The K-means algorithm continues the above iterative steps of Assignment and  
 2825 Update until the centroids  $\boldsymbol{\mu}_i^{(n)}$  ( $i = 1, 2, \dots, K$ ) do not change from those in the  
 2826 previous assignments  $\boldsymbol{\mu}_i^{(n-1)}$  ( $i = 1, 2, \dots, K$ ). When the assignment centroids do  
 2827 not change, we say that the K-means iterative sequence has converged. It can be  
 2828 rigorously proved that this algorithm always converges.

2829 R and Python have K-means functions. Many software packages are freely avail-  
 2830 able for some modified K-means algorithms using different distance formulas. Cli-  
 2831 mate scientists often just use these functions or packages without actually writing  
 2832 an original computer code for a specific K-means algorithm. For example, the R  
 2833 calculation in the previous sub-section used the command `kmeans(mydata, K)` to  
 2834 calculate the cluster results for three points shown in Fig. 3.1.

2835 The K-means algorithm appears very simple, but it was invented only in the  
 2836 1950s. Although the above K-means iterative sequence always converges, the final  
 2837 result may depend on the initial assignment of the centroid:  $C_i^{(1)}$  ( $i = 1, 2, \dots, K$ ),  
 2838 hence each run of the R K-means command `kmeans(mydata, K)` may not always  
 2839 lead to the same solution. Further, a K-means clustering result may not be a global  
 2840 minimum for tWCSS.

2841 Thus, after we have obtained our K-means results, we should try to explain the  
 2842 results and see if they are reasonable. We may also run the code with different  
 2843  $K$  values and check which value  $K$  is the most reasonable. The tWCSS score is a  
 2844 function of  $K$ . We may plot the function tWCSS( $K$ ) and check the shape of the  
 2845 curve. Usually, tWCSS( $K$ ) is a decreasing function, and further,

$$\Delta_K = \text{tWCSS}(K) - \text{tWCSS}(K + 1) > 0, \quad K = 1, 2, \dots \quad (3.20)$$

2846 decreases fast for  $K \leq K_e$ , and then the decrease suddenly slows down or even  
 2847 increases at  $K = K_e$ . In this sense, the curve tWCSS( $K$ ) looks like having an elbow

2848 at  $K_e$ . We would choose this  $K_e$  as the optimal  $K$ . See the elbow in the example  
 2849 of the next sub-section.

2850 Different applications may need different optimization criteria for the selection  
 2851 of the best number of clusters. When you work on your own data, you may try  
 2852 different criteria to choose the best number of clusters. For example, the Silhouette  
 2853 method uses the Silhouette score to measure the similarity level of a point with its  
 2854 own cluster compared to other clusters.

2855 A computer code for the K-means clustering may not always obtain the minimum  
 2856 tWCSS against all partitions. Sometimes, a computer yields a “local” minimum for  
 2857 tWCSS( $K$ ).

2858 In short, when using the K-means method to divide given points into clusters,  
 2859 you should run the K-means code multiple times, check the relevant solution results.  
 2860 Although most K-means code yields a unique solution, multiple solutions sometimes  
 2861 appear, then we need to examine which solution is the most reasonable for our  
 2862 purposes.

### 2863    3.1.3 K-means clustering for the daily Miami weather data

---

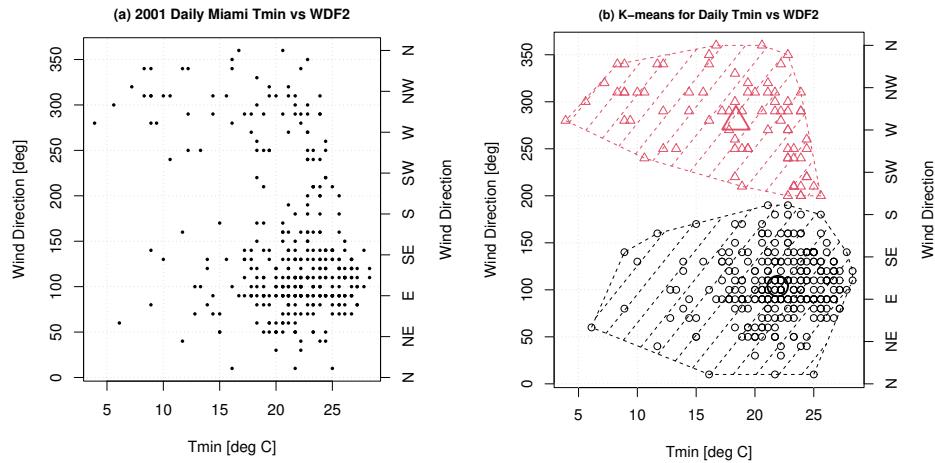
2865 In this sub-section, we present an example using the K-means method to cluster  
 2866 the observed daily weather data at the Miami International Airport, the United  
 2867 States.

#### 2868    3.1.3.1 K-means clustering for $N = 365, K = 2$

2869 Figure 3.2(a) shows a scatter plot of the daily minimum surface air temperature  
 2870 Tmin [in °C] and the direction of the fastest 2-minute wind in a day WDF2 [in  
 2871 degrees] of the Miami International Airport. The wind blowing from the north is  
 2872 zero degree, from the east is 90 degrees, and from the west is 270 degrees. The  
 2873 scatter plot seems to support two clusters or weather regimes: The lower regime  
 2874 showing the wind direction between 0 and 180 degrees, meaning the wind from the  
 2875 east, northeast, or southeast; and another with the wind direction between 180  
 2876 degrees and 360 degrees, meaning that the wind came from the west, northwest, or  
 2877 southwest;

2878 We would like to use the K-means method to make the clustering automatically  
 2879 and justify our above intuitive conclusion from the scatter diagram. This example,  
 2880 although simple, is closer to reality. You can apply the procedures described here  
 2881 to your own data for various clustering purposes, such as defining wildfire zones,  
 2882 ecoregions, climate regimes, and agricultural areas. You can find numerous applica-  
 2883 tion examples by an online search using keywords like “K-means clustering for  
 2884 ecoregions”, or “K-means clustering for climate regimes.”

2885 The daily weather data for the Miami International Airport can be obtained from  
 2886 the Global Historical Climatology Network-Daily (GHCN-D). The station code is  
 2887 USW00012839. The daily data from 2001 to 2020 are included in the book’s master



**Fig. 3.2** (a) Scatter plot of the daily Tmin vs WDF2 for the Miami International Airport in 2001. (b) The K-means clusters of the the daily Tmin vs WDF2 data points when assuming  $K = 2$ .

2888 data file `data.zip` from the book website [www.climatestatistics.org](http://www.climatestatistics.org). The file  
2889 name of the Miami data is

2890 `MiamiIntlAirport2001_2020.csv`

2891 You can access the updated data by an online search for “NOAA Climate Data  
2892 Online USW00012839.” Figure 3.2(b) shows the K-means clustering result for the  
2893 daily Tmin vs WDF2 for the 2001 daily data at the Miami International Airport.  
2894 The two clusters generated by the K-means method agree with our intuition by  
2895 looking at the scatter diagram Fig. 3.2(a). However, the K-means clustering result  
2896 Fig. 3.2(b) provides a more clear picture: A red cluster with each data point indicated  
2897 by a small triangle, and a black cluster with each data point indicated by a  
2898 small circle. The cluster boundaries are linked by dashed lines. The large vertical  
2899 gap at Tmin around 5°C (in winter) between the two clusters indicates that the  
2900 fastest 2-minute wind in a day has a WDF2 angle value around 50°, a northeast  
2901 wind for the red cluster, and around 300°, a northwest wind for the black cluster.

2902 Figure 3.2 may be generated by the following computer codes.

```
2903 #R for Fig. 9.2: K-means clustering for 2001 daily weather
2904 #data at Miami International Airport, Station ID USW00012839
2905 setwd("~/climstats")
2906 dat = read.csv("data/MiamiIntlAirport2001_2020.csv",
2907 header=TRUE)
2908 dim(dat)
2909 #[1] 7305 29
2910 tmin = dat[, 'TMIN']
2911 wdf2 = dat[, 'WDF2']
2912 # plot the scatter diagram Tmin vs WDF2
2913 setEPS() #Plot the data of 150 observations
```

```

2914 postscript("fig0902a.eps", width=5, height=5)
2915 par(mar=c(4.5, 4.5, 2, 4.5))
2916 plot(tmin[2:366], wdf2[2:366],
2917       pch =16, cex = 0.5,
2918       xlab = 'Tmin[degC]',
2919       ylab = 'WindDirection[deg]', grid())
2920 title('(a) 2001 Daily Miami Tmin vs WDF2',
2921       cex.main = 0.9, line = 1)
2922 axis(4, at = c(0, 45, 90, 135, 180, 225, 270, 315, 360),
2923       lab = c('N', 'NE', 'E', 'SE', 'S', 'SW', 'W', 'NW', 'N'))
2924 mtext('WindDirection', side = 4, line =3)
2925 dev.off()
2926 #K-means clustering
2927 K = 2 #assuming K = 2, i.e., 2 clusters
2928 mydata = cbind(tmin[2:366], wdf2[2:366])
2929 fit = kmeans(mydata, K) # K-means clustering
2930 #Output the coordinates of the cluster centers
2931 fit$centers
2932 #1 18.38608 278.8608
2933 #2 21.93357 103.9161
2934 fit$tot.withinss # total WCSS
2935 #[1] 457844.9 for # the value may vary for each run
2936
2937 #Visualize the clusters by kmeans.ani()
2938 mycluster <- data.frame(mydata, fit$cluster)
2939 names(mycluster)<-c('Tmin[degC]',
2940                      'WindDirection[deg]',
2941                      'Cluster')
2942 library(animation)
2943 par(mar = c(4.5, 4.5, 2, 4.5))
2944 kmeans.ani(mycluster, centers = K, pch=1:K, col=1:K,
2945             hints = '')
2946 title(main=
2947       "(b) K-means Clusters for Daily Tmin vs WDF2",
2948       cex.main = 0.8)
2949 axis(4, at = c(0, 45, 90, 135, 180, 225, 270, 315, 360),
2950       lab = c('N', 'NE', 'E', 'SE', 'S', 'SW', 'W', 'NW', 'N'))
2951 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') #centers 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 [°C]", labelpad = 10)
plt.ylabel(r"Wind Direction [°]", 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()

```

### 2953 3.1.3.2 Why $K = 2$ for the Miami weather data?

2954 The K-means computer code for the Miami Tmin and WDF2 data shows two  
 2955 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  
 2956 total WCSS is 457844.9. To justify our choice of  $K = 2$ , we compute tWCSS, the  
 2957 total WCSS, for different numbers of  $K$ , and then use the elbow rule to determine  
 2958 the best  $K$  value. Figure 3.3(a) shows the tWCSS scores for  $K = 2, 3, \dots, 8$ . We  
 2959 can see that the elbow appears at  $K = 2$ , since total WCSS decrement suddenly  
 2960 slows down from  $K = 2$ . Thus, the elbow rule has justified our choice of  $K = 2$ .

2961 Another method to choose  $K$  is the knee rule, which is determined by the vari-  
 2962 ances explained by  $K$  clusters. When  $K = 1$ , tWCSS is the spatial variance of the  
 2963 entire data. When  $K = N$ , tWCSS = 0, where  $N$  is the total number of the data  
 2964 points. Usually, we have

$$\text{tWCSS}[1] \geq \text{tWCSS}[K] \geq \text{tWCSS}[N] \quad (3.21)$$

2965 for any  $1 \leq K \leq N$ . Then,

$$\text{pWCSS}[K] = \frac{\text{tWCSS}[1] - \text{tWCSS}[K]}{\text{tWCSS}[1]} \times 100\% \quad (3.22)$$

2966 is defined as the percentage of variance explained by the  $K$  clusters. Usually,  
 2967  $\text{pWCSS}[K]$  is an increasing function of  $K$ . When the increase rate suddenly slows  
 2968 down from  $K$ , a shape like a knee appears in the  $\text{pWCSS}[K]$  curve. We then use  
 2969 this  $K$  as the optimal number of centers for our K-means clustering. This is the  
 2970 so called the knee rule. Figure 3.3(b) shows that the increase of  $\text{pWCSS}[K]$  dra-  
 2971 matically slows down at  $K = 2$ . Thus, we choose  $K = 2$ , which provides another  
 2972 justification why  $K = 2$ .

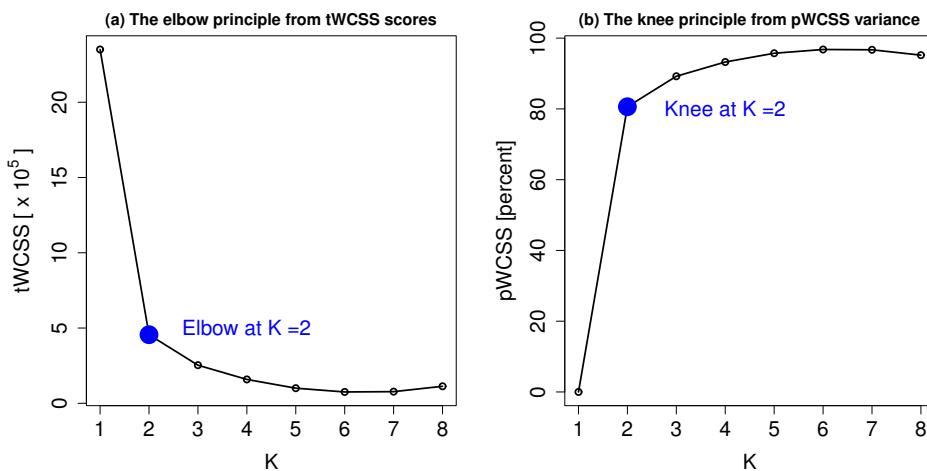


Fig. 3.3 (a) tWCSS scores for different  $K$ . (b) pWCSS variances for different  $K$ .

2973 Figure 3.3 may be generated by the following computer codes.

```
2974 #R plot Fig. 9.3: tWCSS(K) and pWCSS(K)
2975 twcss = c()
2976 for(K in 1:8){
2977   mydata=cbind(tmax[2:366], wdf2[2:366])
2978   twcss[K] = kmeans(mydata, K)$tot.withinss
2979 }
2980 twcss
2981 par(mar = c(4.5, 6, 2, 0.5))
2982 par(mfrow=c(1,2))
2983 plot(twcss/100000, type = 'o', lwd = 2,
2984       xlab = 'K', ylab = bquote('tWCSS[~x, ~ 10^5 ~]'),
2985       main = '(a) The elbow principle from tWCSS scores',
2986       cex.lab = 1.5, cex.axis = 1.5)
2987 points(2, twcss[2]/100000, pch =16, cex = 3, col = 'blue')
2988 text(4, 5, 'Elbow at K=2', cex = 1.5, col = 'blue')
2989 #compute percentage of variance explained
2990 pWCSS = 100*(twcss[1]- twcss)/twcss[1]
2991 plot(pWCSS, type = 'o', lwd = 2,
2992       xlab = 'K', ylab = 'pWCSS[percent]',
2993       main = '(b) The knee principle from pWCSS variance',
2994       cex.lab = 1.5, cex.axis = 1.5)
2995 points(2, pWCSS[2], pch =16, cex = 3, col = 'blue')
2996 text(4, 80, 'Knee at K=2', cex = 1.5, col = 'blue')
2997 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()

```

2998

2999

### 3.1.3.3 Steps of the K-means clustering data analysis

3000 From the previous sub-sections, we may summarize that the K-means method for  
 3001 data analysis consists of the following steps.

- 3002 (i) Data preparation: Write your data into an  $N \times 2$  matrix, where  $N$  is the total  
 3003 number of data points, and 2 is for two variables, i.e., in a 2-dimensional

space. The real value matrix must have no missing data. Although our previous examples are for 2-dimensional data, the K-means clustering can be used for higher dimensional data. On the Internet, you can easily find examples of the K-means for data of three or more variables.

- (ii) Preliminary data analysis: You may wish to make some preliminary simple analyses before applying the K-means method. These analyses may include plotting a scatter diagram, plot a histogram for each variable, and compute mean, standard deviation and quantiles for each variable. If the data for the two variables have different units, you may convert the data into their standardized anomalies. The standardized anomaly computing procedure is called “scale” in the machine learning community. For example, the R command `scale(x)` yields the standard anomalies of data sequence `x`, i.e, divided by the standard deviation after a subtraction of mean. Mathematically speaking, the scaling procedure is preferred when two variables have different units, because the  $\text{unit1}^2 + \text{unit2}^2$  in tWCSS usually do not have a good interpretation. In contrast, standardized anomalies are dimensionless and have no units. Non-dimensional data are often preferred when applying statistical methods. However, our goal of the K-means clustering is to find patterns without labeling the data, which is an unsupervised learning. The K-means results from scaled data and unscaled data may not be the same. As long as the K-means result is reasonable and interpretable, regardless of the scaled or unscaled data, we adopt the result. For example, in our Miami Tmin and WDF2 example in the previous sub-section, we found the K-means clusters for both scaled and unscaled data, but we have chosen the result from the unscaled data from our view point of interpretation. Therefore, we keep in mind that the K-means is a result-oriented method.
- (iii) Test K-means clustering: Apply a computer code of K-means to the prepared  $N \times 2$  data matrix with a  $K$  value inspired from the scatter diagram or science background of the data. You may apply the code once to the scaled data and another to the unscaled data, and see if the results make sense. Run your code a few times and see if the results vary.
- (iv) Determine the optimal  $K$  by the elbow rule: Compute the tWCSS scores for different  $K$  and find the elbow location of the  $\text{tWCSS}(K)$  curve. Apply the K-means computer code for this optimal  $K$  value. Compare this  $K$  with the  $K$  in the previous step. Check the details of the clustering results, such as cluster assignment, tWCSS, centers.
- (v) Visualize the K-means clustering results: Use a computer code to plot the  $K$  clusters. Different visualization codes are available for your choice.
- (vi) Interpret the final result: Identify the characteristics of the clusters and interpret them with text or more figures.

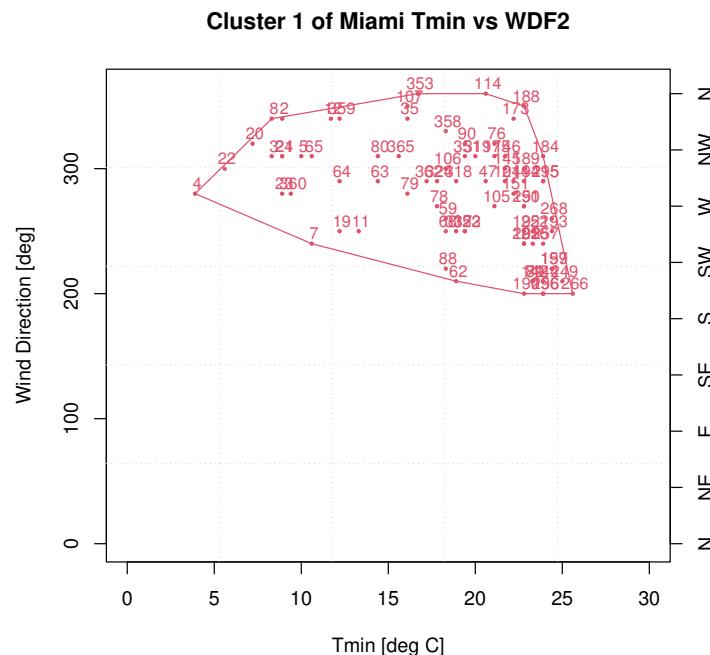
Applying the K-means computer code is very easy, but in the K-means data analysis, we should avoid the following common mistakes:

- Skipped the preliminary data analysis step: Some people may quickly apply the

3047 K-means computer code and treat the K-means results as unique, hence regard  
 3048 their plot of the clusters as the end of the analysis.

- 3049 • Missed the interpretation: Interpretation step is needed, short or long. This is  
 3050 our purpose anyway.  
 3051 • Forgot to scale the data: We should always analyze the scaled data. At the same  
 3052 time, because of the interpretation requirement, we should also check the re-  
 3053 sults from the unscaled data, and see which result is more reasonable.  
 3054 • Overlooked the justification of the optimal  $K$ : A complete K-means clustering  
 3055 should have this step, which helps us to explore all the possible clusters.

3056 In addition to using the standard K-means computer code you can find from the  
 3057 Internet or that built in R or Python, you sometimes may wish to make special  
 3058 plots for your clustering results with some specific features, such as a single cluster  
 3059 on a figure with the marked data order, as shown in Fig. 3.4. Convex hull is a  
 3060 simple method to plot a cluster. By definition, the convex hull of a set of points is  
 3061 the smallest convex polygon that encloses all of the points in the set. Below is the  
 3062 a computer code to plot Fig. 3.4 using the convex hull method.



3063     Figure 3.4 may be generated by the following computer codes.

```

3064 #R plot Fig. 9.4: Convex hull for a cluster
3065 setwd("~/climstats")
3066 dat = read.csv("data/MiamiIntlAirport2001_2020.csv",
3067                 header=TRUE)
3068 dim(dat)
3069 #[1] 7305   29
3070 tmin = dat[, 'TMIN']
3071 wdf2 = dat[, 'WDF2']
3072 #K-means clustering
3073 K = 2 #assuming K = 2, i.e., 2 clusters
3074 mydata = cbind(tmin[2:366], wdf2[2:366]) #2001 data
3075 clusterK = kmeans(mydata, K) # K-means clustering
3076 mycluster <- data.frame(mydata, clusterK$cluster)
3077 plotdat = cbind(1:N, mycluster)
3078
3079 par(mar = c(4.5, 6, 2, 0.5)) #set up the plot margin
3080 i = 1 # plot Cluster 1
3081 N = 365 #Number of data points
3082 X = plotdat[which(mycluster[,3] == i), 1:3]
3083 colnames(X)<-c('Day', 'Tmin[degC]', 'WDF2[deg]')
3084 plot(X[,2:3], pch = 16, cex = 0.5, col = i,
3085       xlim = c(0, 30), ylim = c(0, 365),
3086       main = 'Cluster 1 of Miami Tmin vs WDF2')
3087 grid(5, 5)
3088 #chull() finds the boundary points of a convex hull
3089 hpts = chull(X[, 2:3])
3090 hpts1 = c(hpts, hpts[1]) #close the boundary
3091 lines(X[hpts1, 2:3], col = i)
3092 for(j in 1:length(X[,1])){
3093   text(X[j,2], X[j,3] + 8, paste("", X[j,1]),
3094         col = i, cex = 0.8)
3095 } #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])

```

3096

3097     The K-means method can not only be used for weather data clustering and  
 3098     climate classification, but also for many other purposes, such as weather forecasting,  
 3099     storm identification, and more. In many applications, K-means is part of a machine  
 3100     learning algorithm, and takes a modified version, such as the incremental K-means  
 3101     for weather forecasting and the weighted K-means for climate analysis.

3102

## 3.2 Support vector machine

3103

3104     The K-means clustering can organize a suite of given data points into  $K$  clusters.  
 3105     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.

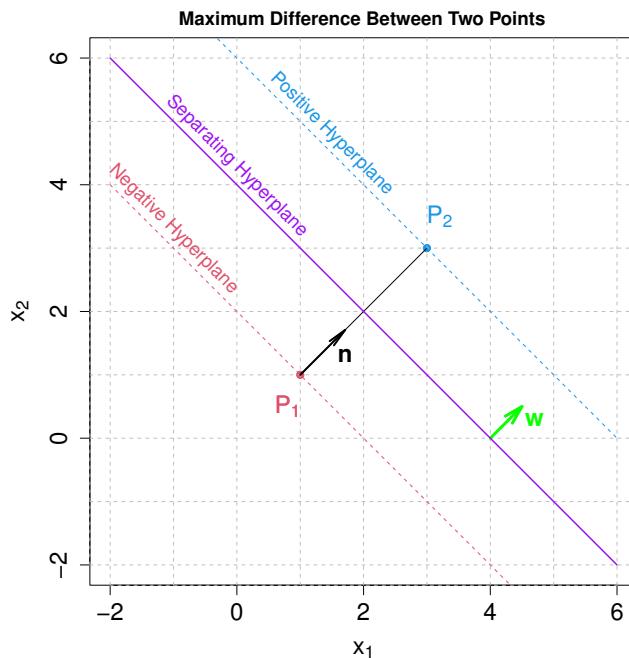
In the following we will illustrate these SVM learning ideas using three examples: (i) A trivial case of two data points, (ii) training and prediction for a system of three data points, and (iii) a more general case of many data points.

### 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

$$x_1 + x_2 = 4, \quad (3.23)$$

that for the positive hyperplane (i.e., the blue dashed line) is

$$x_1 + x_2 = 6, \quad (3.24)$$

<sup>3165</sup> and that for the negative hyperplane (i.e., the red dashed line) is

$$x_1 + x_2 = 2. \quad (3.25)$$

<sup>3166</sup> The unit vector in the perpendicular bisector direction from  $P_1$  to  $P_2$  is

$$\mathbf{n} = \left( \frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right). \quad (3.26)$$

<sup>3167</sup> This vector is a normal vector for the separating hyperplane. The equation of the  
<sup>3168</sup> separating hyperplane may be expressed in the normal vector form

$$\mathbf{n} \cdot \mathbf{x} = 2\sqrt{2}. \quad (3.27)$$

<sup>3169</sup> This can further be written in the following form

$$\mathbf{w} \cdot \mathbf{x} - b = 0, \quad (3.28)$$

<sup>3170</sup> where

$$\mathbf{w} = (0.5, 0.5) \quad (3.29)$$

<sup>3171</sup> and

$$b = 2. \quad (3.30)$$

<sup>3172</sup> An advantage of this form is that the equations for the positive and negative hy-  
<sup>3173</sup> perplanes can be simply written in the form

$$\mathbf{w} \cdot \mathbf{x} - b = \pm 1, \quad (3.31)$$

<sup>3174</sup> and that for the separating hyperplane can be written as

$$\mathbf{w} \cdot \mathbf{x} - b = 0. \quad (3.32)$$

<sup>3175</sup> In the above,  $\mathbf{w}$  is called a normal vector and is related with the unit normal vector  
<sup>3176</sup>  $\mathbf{n}$  as follows:

$$\mathbf{w} = -\frac{1}{\sqrt{2}} \mathbf{n}. \quad (3.33)$$

<sup>3177</sup> The distance between the positive and negative hyperplanes is

$$D_m = \frac{2}{|\mathbf{w}|} = 2\sqrt{2}, \quad (3.34)$$

<sup>3178</sup> is the margin size. The maximum difference between  $P_1$  and  $P_2$  is quantified by  
<sup>3179</sup>  $D_m$ . The primary goal of SVM is to maximize  $D_m$  when training data are given.

<sup>3180</sup> The distance of the origin to the separating hyperplane is

$$D_o = -\frac{b}{|\mathbf{w}|} = 2\sqrt{2}. \quad (3.35)$$

<sup>3181</sup> These formulations of the hyperplanes and the quantities  $\mathbf{w}, b, D_m$  can be used  
<sup>3182</sup> for the general SVM mathematical formulation, which are described in the next  
<sup>3183</sup> sub-section.

<sup>3184</sup> Figure 3.5 may be generated by the following computer code.

```

3185 #R plot Fig. 9.5: Maximum difference between two points
3186 x = matrix(c(1, 1, 3, 3),
3187   ncol = 2, byrow = TRUE)#Two points
3188 y= c(-1, 1) #Two labels -1 and 1
3189 #Plot the figure and save it as a .eps file
3190 setEPS()
3191 postscript("fig0905.eps", height=7, width=7)
3192 par(mar = c(4.5, 4.5, 2.0, 2.0))
3193 plot(x, col = y + 3, pch = 19,
3194   xlim = c(-2, 6), ylim = c(-2, 6),
3195   xlab = bquote(x[1]), ylab = bquote(x[2]),
3196   cex.lab = 1.5, cex.axis = 1.5,
3197   main = "Maximum\u00d7Difference\u00d7Between\u00d7Two\u00d7Points")
3198 axis(2, at = (-2):6, tck = 1, lty = 2,
3199   col = "grey", labels = NA)
3200 axis(1, at = (-2):6, tck = 1, lty = 2,
3201   col = "grey", labels = NA)
3202 segments(1, 1, 3, 3)
3203 arrows(1, 1, 1.71, 1.71, lwd = 2,
3204   angle = 9, length= 0.2)
3205 text(1.71, 1.71-0.4, quote(bold('n')), cex = 1.5)
3206 arrows(4, 0, 4 + 0.5, 0 + 0.5, lwd = 3,
3207   angle = 15, length= 0.2, col = 'green' )
3208 text(4.7, 0.5 -0.2, quote(bold('w')),
3209   cex = 1.5, col = 'green')
3210 x1 = seq(-2, 6, len = 31)
3211 x20 = 4 - x1
3212 lines(x1, x20, lwd = 1.5, col = 'purple')
3213 x2m = 2 - x1
3214 lines(x1, x2m, lty = 2, col = 2)
3215 x2p = 6 - x1
3216 lines(x1, x2p, lty = 2, col = 4)
3217 text(1-0.2,1-0.5, bquote(P[1]), cex = 1.5, col = 2)
3218 text(3+0.2,3+0.5, bquote(P[2]), cex = 1.5, col = 4)
3219 text(1-0.2,1-0.5, bquote(P[1]), cex = 1.5, col = 2)
3220 text(0,4.3, 'Separating\u00d7Hyperplane',
3221   srt = -45, cex = 1.2, col = 'purple')
3222 text(1.5, 4.8, 'Positive\u00d7Hyperplane',
3223   srt = -45, cex = 1.2, col = 4)
3224 text(-1, 3.3, 'Negative\u00d7Hyperplane',
3225   srt = -45, cex = 1.2, col = 2)
3226 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[0], color = 'r')
plt.scatter(x[1],x[1], 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()

```

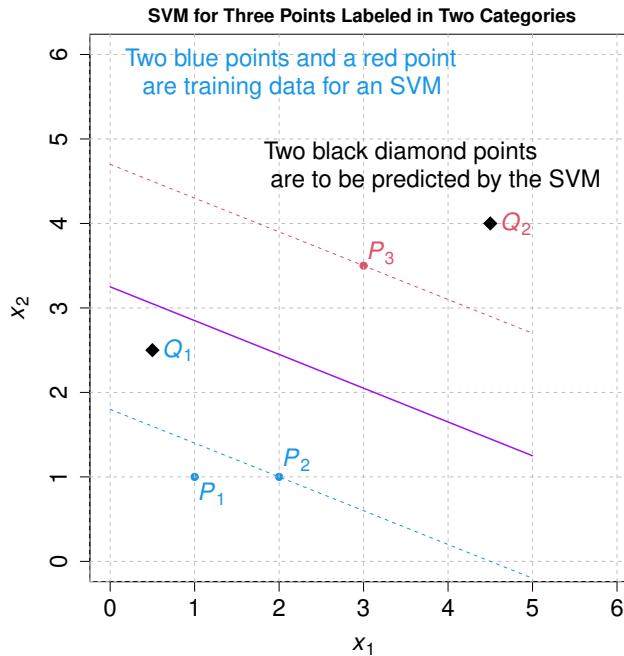
3227

### 3.2.1 SVM for a system of three points labeled in two categories

3228 Let us progress from a two-point system to a three-point system and build an  
 3229 SVM using a computer code. Figure 3.6 shows the three points labeled in two  
 3230 categories. The two blue points are  $P_1(1,1), P_2(2,1)$  which are in the first category  
 3231 corresponding to the categorical value  $y = 1$ . The red point  $P_3(3,3.5)$  is in the  
 3232 second category corresponding to  $y = -1$ . Training these data is to maximize the  
 3233 margin size  $D_m$ , which results in the following SVM parameters.  
 3234

$$w = (-0.28 - 0.69), \quad b = 2.24, \quad D_m = 2.69. \quad (3.36)$$

3236 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.

3237 We acquire two new data points  $Q_1(0.5, 2.5)$  and  $Q_2(4.5, 4.0)$ . We wish to use  
3238 the trained SVM to predict whether  $Q_1(0.5, 2.5)$  corresponds to 1 (i.e., the first  
3239 category) or -1 (i.e., the second category). The SVM prediction results are shown  
3240 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  
3241 category. Geometrically, this prediction result in Fig. 3.6 is obvious and has nothing  
3242 surprising. The significance of this example is its systematic use of the SVM method  
3243 and its computer algorithm. Although this system is simple and has only three  
3244 points, solving this problem by hand is very difficult. The following computer code  
3245 solves this problem and generates Fig. 3.6.

```

3246 #R plot Fig. 9.6: SVM for three points
3247 #training data x
3248 x = matrix(c(1, 1, 2, 1, 3, 3.5),
3249   ncol = 2, byrow = TRUE)
3250 y = c(1, 1, -1) #two categories 1 and -1
3251 plot(x, col = y + 3, pch = 19,
3252   xlim = c(-2, 8), ylim = c(-2, 8))
3253 library(e1071)
3254 dat = data.frame(x, y = as.factor(y))
3255 svm3P = svm(y ~ ., data = dat,
3256   kernel = "linear", cost = 10,
3257   scale = FALSE,
3258   type = 'C-classification')
3259 svm3P #This is the trained SVM

```

```

3260 xnew = matrix(c(0.5, 2.5, 4.5, 4),
3261           ncol = 2, byrow = TRUE) #New data
3262 predict(svm3P, xnew) #prediction using the trained SVM
3263 # 1 2
3264 # 1 -1
3265
3266 # Find hyperplane, normal vector, and SV (wx + b = 0)
3267 w = t(svm3P$coeffs) %*% svm3P$SV
3268 w
3269 #[1,] -0.2758621 -0.6896552
3270 b = svm3P$rho
3271 b
3272 #[1] -2.241379
3273 2/norm(w, type = '2') #maximum margin of separation
3274 #[1] 2.692582
3275
3276 x1 = seq(0, 5, len = 31)
3277 x2 = (b - w[1]*x1)/w[2]
3278 x2p = (1 + b - w[1]*x1)/w[2]
3279 x2m = (-1 + b - w[1]*x1)/w[2]
3280 x20 = (b - w[1]*x1)/w[2]
3281 #plot the SVM results
3282 setEPS()
3283 postscript("fig0906.eps", height=7, width=7)
3284 par(mar = c(4.5, 4.5, 2.0, 2.0))
3285 plot(x, col = y + 3, pch = 19,
3286       xlim = c(0, 6), ylim = c(0, 6),
3287       xlab = bquote(x[1]), ylab = bquote(x[2]),
3288       cex.lab = 1.5, cex.axis = 1.5,
3289       main = 'SVM for three points labeled in two categories')
3290 axis(2, at = (-2):8, tck = 1, lty = 2,
3291       col = "grey", labels = NA)
3292 axis(1, at = (-2):8, tck = 1, lty = 2,
3293       col = "grey", labels = NA)
3294 lines(x1, x2p, lty = 2, col = 4)
3295 lines(x1, x2m, lty = 2, col = 2)
3296 lines(x1, x20, lwd = 1.5, col = 'purple')
3297 xnew = matrix(c(0.5, 2.5, 4.5, 4),
3298               ncol = 2, byrow = TRUE)
3299 points(xnew, pch = 18, cex = 2)
3300 for(i in 1:2){
3301   text(xnew[i,1] + 0.5, xnew[i,2] , paste('Q',i),
3302         cex = 1.5, col = 6-2*i)
3303 }
3304 text(2.2,5.8, "Two blue points and a red point
3305 are training data for an SVM",
3306       cex = 1.5, col = 4)
3307 text(3.5,4.7, "Two black diamond points
3308 are to be predicted by the SVM",
3309       cex = 1.5)
3310 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()

```

---

3312    **3.2.2 SVM mathematical formulation for a system of many**  
 3313    **points in two categories**

---

3315    The previous quantification of maximum separation between two points or three  
 3316    points can be generalized to the separation of many points of two categories labeled  
 3317     $(1, -1)$ , or  $((1, 2)$ , or  $(A, B)$ , or another way. What is a systematic formulation of  
 3318    mathematics to make the best separation, by maximizing the margin? How can we  
 3319    use a computer code to implement the separation? This sub-section attempts to  
 3320    answer these questions.

3321    We have  $n$  training data points:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), \quad (3.37)$$

3322    where  $x_i$  is a  $p$ -dimensional vector (i.e., using  $p$  parameters to describe every data  
 3323    point in a category), and  $y_i$  is equal to 1 or  $-1$ , a category indicator,  $i = 1, 2, \dots, n$ .  
 3324    The separating hyperplane has its linear equation as follows

$$\mathbf{w} \cdot \mathbf{x} - b = 0. \quad (3.38)$$

3325    Our SVM algorithm is to find the  $p$ -dimensional normal vector  $\mathbf{w}$  and a real scalar  
 3326     $b$  so that the distance

$$D_m = \frac{2}{|\mathbf{w}|} \quad (3.39)$$

3327    between the positive and negative hyperplanes

$$\mathbf{w} \cdot \mathbf{x} - b = \pm 1 \quad (3.40)$$

3328    is maximized.

3329    The maximization is the domain in the  $p$ -dimensional space below the negative  
 3330    hyperplane and above the positive hyperplane, i.e.,

$$\mathbf{w} \cdot \mathbf{x}_i - b \geq 1 \text{ when } y_i = 1 \quad (3.41)$$

3331    and

$$\mathbf{w} \cdot \mathbf{x}_i - b \leq -1 \text{ when } y_i = -1. \quad (3.42)$$

3332    This definition of the maximization implies that no training data points are located  
 3333    between the positive and negative hyperplanes. Thus, the solution of the maximization  
 3334    of  $D_m = 2/|\mathbf{w}|$  must occur at the boundary of the domain, i.e., the positive  
 3335    and negative hyperplanes. The points on these hyperplanes are called the support  
 3336    vectors (SV). R or Python SVM code can find these SVs.

3337    Figure 3.7 shows an example of an SVM training with 20 training data points,  
 3338    and an SVM prediction for 3 data points. The data are given below:

3339	X1	X2	y	
3340	1	1	6.0	1
3341	2	2	8.0	1
3342	3	3	7.5	1

```

3343 4 1 8.0 1
3344 5 4 9.0 1
3345 6 5 9.0 1
3346 7 3 7.0 1
3347 8 5 9.0 1
3348 9 1 5.0 1
3349 10 5 3.0 2
3350 11 6 4.0 2
3351 12 7 4.0 2
3352 13 8 6.0 2
3353 14 9 5.0 2
3354 15 10 6.0 2
3355 16 5 0.0 2
3356 17 6 5.0 2
3357 18 8 2.0 2
3358 19 2 2.0 2
3359 20 1 1.0 2

```

3360 The new data to be predicted are

```

3361 [1,] 0.5 2.5
3362 [2,] 7.0 2.0
3363 [3,] 6.0 9.0

```

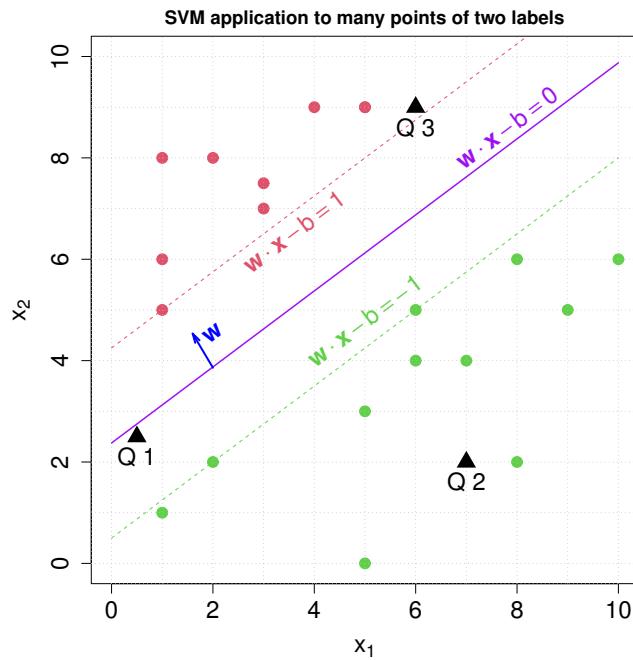
3364 The following computer code can do the SVM calculation and plot Fig. 3.7.

3365 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  
 3366 wish to use the trained SVM to find out which point belongs to what category. The  
 3367 SVM prediction result is shown in Fig. 3.7:  $Q_1(0.5, 2.5)$  and  $Q_2(7.5, 2.0)$  belong to  
 3368 Category 1, and  $Q_3(6.0, 9.0)$  to Category 2. Again, this prediction result is obvious  
 3369 to us when the new data points are plotted, our eyes can see the points, and  
 3370 our brains can make the decisions based on our trained intelligence. The SVM  
 3371 applications in the real world often involve data of higher dimensions, in which  
 3372 case the visual decision is impossible, and requires us to predict the category for  
 3373 the new data without plotting. The prediction is based on the training data, new  
 3374 data, and SVM algorithm without visualization. For example, the R prediction of  
 3375 this problem can be done by the command `predict(svmP, xnew)` where `xnew` is  
 3376 the new data in a  $3 \times 2$  matrix, and `svmP` is the trained SVM. This command is  
 3377 included in the computer code for generating Fig. 3.7.

```

3378 #R plot Fig. 9.7: SVM for many points
3379 #Training data x and y
3380 x = matrix(c(1, 6, 2, 8, 3, 7.5, 1, 8, 4, 9, 5, 9,
3381           3, 7, 5, 9, 1, 5,
3382           5, 3, 6, 4, 7, 4, 8, 6, 9, 5, 10, 6,
3383           5, 0, 6, 5, 8, 2, 2, 2, 1, 1),
3384           ncol = 2, byrow = TRUE)
3385 y= c(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.

```

3386      1, 1, 1,
3387      2, 2, 2, 2, 2, 2,
3388      2, 2, 2, 2, 2)
3389
3390 library(e1071)
3391 dat = data.frame(x, y = as.factor(y))
3392 svmP = svm(y ~ ., data = dat,
3393             kernel = "linear", cost = 10,
3394             scale = FALSE,
3395             type = 'C-classification')
3396 svmP
3397 #Number of Support Vectors: 3
3398 svmP$SV #SVs are #x[9,], x[17,], x[19,]
3399 #9 1 5
3400 #17 6 5
3401 #19 2 2
3402
3403 # Find SVM parameters: w, b, SV (wx+c=0)
3404 w = t(svmP$coefs) %*% svmP$SV
3405 # In essence this finds the hyper plane that separates our points
3406 w
3407 #[1,] -0.39996 0.53328
3408 b = svmP$rho
3409 b
3410 #[1] 1.266573
3411 2/norm(w, type ='2')

```

```

3412 #[1] 3.0003 is the maximum margin
3413 x1 = seq(0, 10, len = 31)
3414 x2 = (b - w[1]*x1)/w[2]
3415 x2p = (1 + b - w[1]*x1)/w[2]
3416 x2m = (-1 + b - w[1]*x1)/w[2]
3417 x20 = (b - w[1]*x1)/w[2]
3418
3419 #plot the svm results
3420 setEPS()
3421 postscript("fig0907.eps", height=7, width=7)
3422 par(mar = c(4.5, 4.5, 2.0, 2.0))
3423 plot(x, col = y + 9, pch = 19, cex = 1.5,
3424       xlim = c(0, 10), ylim = c(0, 10),
3425       xlab = bquote(x[1]), ylab = bquote(x[2]),
3426       cex.lab = 1.5, cex.axis = 1.5,
3427       main = 'SVM application to many points of two labels')
3428 axis(2, at = 0:10, tck = 1, lty = 3,
3429       col = "grey", labels = NA)
3430 axis(1, at = 0:10, tck = 1, lty = 3,
3431       col = "grey", labels = NA)
3432 lines(x1, x2p, lty = 2, col = 10)
3433 lines(x1, x2m, lty = 2, col = 11)
3434 lines(x1, x20, lwd = 1.5, col = 'purple')
3435 thetasvm = atan(-w[1]/w[2])*180/pi
3436 thetasvm
3437 #[1] 36.8699 #36.9 deg angle of the hyperplane
3438 #linear equations for the hyperplanes
3439 delx = 1.4
3440 dely = delx * (-w[1]/w[2])
3441 text(5 + 2*delx, 6.5 + 2*dely, bquote(bold(w%.%x) - b == 0),
3442       srt = thetasvm, cex = 1.5, col = 'purple')
3443 text(5 - delx, 7.6 - dely, bquote(bold(w%.%x) - b == 1),
3444       srt = thetasvm, cex = 1.5, col = 10)
3445 text(5, 4.8, bquote(bold(w%.%x) - b == -1),
3446       srt = thetasvm, cex = 1.5, col = 11)
3447 #normal direction of the hyperplanes
3448 arrows(2, 3.86, 2 + w[1], 4 + w[2], lwd = 2,
3449         angle = 15, length= 0.1, col = 'blue' )
3450 text(2 + w[1] + 0.4, 4 + w[2], bquote(bold(w)),
3451       srt = thetasvm, cex = 1.5, col = 'blue')
3452
3453 #new data points to be predicted
3454 xnew = matrix(c(0.5, 2.5, 7, 2, 6, 9),
3455                 ncol = 2, byrow = TRUE)
3456 points(xnew, pch = 17, cex = 2)
3457 predict(svmP, xnew) #Prediction
3458 #1 2 3
3459 #2 2 1
3460
3461 for(i in 1:3){
3462   text(xnew[i,1], xnew[i,2] - 0.4 ,
3463         paste('Q',i), cex = 1.5)
3464 }
3465 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 = [1,1,1,1,1,1,1,1,2,2,2,2,2,2,2,2,2,2,2,2]
# 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()

```

3467

SVM maximizes  $2/\|\mathbf{w}\|$ , i.e., minimizes  $\|\mathbf{w}\|$ . The above computer code shows that  $\mathbf{w} = (-0.39996, 0.53328)$ . This optimization is reached at the three SVs  $P_9(1, 5)$ ,  $P_{17}(6, 5)$  and  $P_{19}(2, 2)$ . We can verify that the end points of these three SVs are on the two hyperplanes, by plugging in these coordinates into the following two linear equations

$$\mathbf{w} \cdot \mathbf{x} - b = \pm 1. \quad (3.43)$$

3468 Reversely, these equations have three parameters  $w_1, w_2, b$ , which can be determined  
 3469 by three support vectors (SVs).

3470 The linear SVM separation can have nonlinear extensions, such as circular or  
 3471 spherical hypersurfaces. Nonlinear SVM is apparently needed when the data cannot  
 3472

3477 be easily separated by linear hyperplanes. Both R and Python codes have function  
3478 parameters for the nonlinear SVM. Another SVM extension is from two classes to  
3479 many classes. The multi-class SVM and nonlinear SVM are beyond the scope of  
3480 this book. The interested readers are referred to the designated machine learning  
3481 books, such as those listed at the section of References and Further Readings at  
3482 the end of this chapter.

### 3483 **3.3 Random forest method for classification and** 3484 **regression**

---

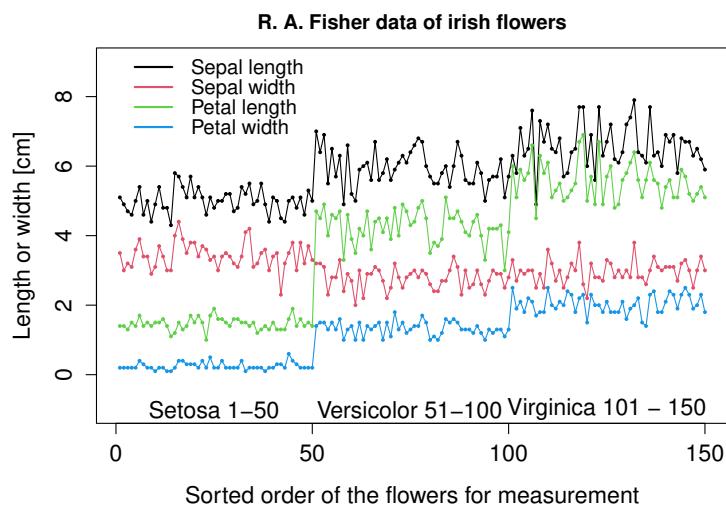
3485 Random forest (RF) is a popular machine learning algorithm and belongs to the  
3486 class of supervised learning. It uses many decision trees trained from the given data,  
3487 called the training data. The training data usually include categorical data, such  
3488 as weather types (e.g., rainy, cloudy, and sunny) as labels, and numeric data, such  
3489 as historical instrument observations (e.g., atmospheric pressure, humidity, wind  
3490 speed, wind direction, and air temperature). These data form many logical decision  
3491 trees that decide under what sets of conditions the weather is rainy, cloudy, or  
3492 sunny. Many decision trees form a forest. Multiple decision nodes and branches for  
3493 each tree are determined and trained by using a set of random sampling procedures.  
3494 The random sampling and the decision trees lead to the name of random forest.  
3495 The training step results in an RF model which are a set of decision trees. In the  
3496 prediction step, you submit your new data to the the trained RF model whose each  
3497 decision tree makes a vote for a category based on your new data. If the category  
3498 “tomorrow is rainy” receives the most votes, then you have predicted a rainy day  
3499 tomorrow. In this sense, RF is an ensemble learning method. Because of this nature  
3500 of ensemble predictions, an RF algorithm should ensure the votes among the trained  
3501 trees should have as little correlation as possible.

3502 The above is a layman’s description of RF for classification, when the RF objective  
3503 is for categorical data, or called factor in data types. RF can also be used to fill  
3504 in the missing numeric real values. The missing values, or missing data in climate  
3505 science, are often denoted by NA, NaN, or -99999. RF makes a prediction for the  
3506 missing data. The prediction is an ensemble mean from the trained RF trees. This  
3507 is an RF regression, which produces numerical results.

#### 3509 **3.3.1 RF flower classification for a benchmark iris dataset**

---

3510 We use the popular iris flower dataset (Fisher 1936) frequently used for the RF  
3511 teaching as an example to explain the RF computing procedures and to interpret  
3512 the RF results. Sir Ronald A. Fisher (1890 – 1962) was a British statistician and  
3513 biologist. The Fisher dataset contains three iris species: *setosa*, *versicolor*, *virginica*.  
3514 The data are the length and width of sepal and petal of a flower. The first two rows  
3515 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).

```

3517 #   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
3518 #1      5.1        3.5       1.4        0.2  setosa
3519 #2      4.9        3.0       1.4        0.2  setosa

```

3520 Figure 3.8 shows the entire dataset.

3521 Figure 3.8 may be generated by the following computer code.

```

3522 #R plot Fig. 9.8: R.A. Fisher data of three iris species
3523 setwd('~/climstats')
3524 data(iris) #read the data already embedded in R
3525 dim(iris)
3526 #[1] 150  5
3527 iris[1:2,] # Check the first two rows of the data
3528 #   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
3529 #1      5.1        3.5       1.4        0.2  setosa
3530 #2      4.9        3.0       1.4        0.2  setosa
3531
3532 str(iris) # Check the structure of the data
3533 #'data.frame': 150 obs. of 5 variables:
3534 #\$ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
3535 #\$ Species     : Factor w/ 3 levels "setosa","versicolor",...
3536
3537 setEPS() #Plot the data of 150 observations
3538 postscript("fig0908.eps", width=7, height=5)
3539 par(mar= c(4.5, 4.5, 2.5, 0.2))
3540 plot(iris[,1], type = 'o', pch = 16, cex = 0.5, ylim = c(-1, 9),
3541       xlab = 'Sorted order of the flowers for measurement',
3542       ylab = 'Length or width [cm]',
3543       main = 'R.A.Fisher data of irish flowers', col = 1,
3544       cex.lab = 1.3, cex.axis = 1.3)

```

```

3545 lines(iris[,2], type = 'o', pch = 16, cex = 0.5, col=2)
3546 lines(iris[,3], type = 'o', pch = 16, cex = 0.5, col=3)
3547 lines(iris[,4], type = 'o', pch = 16, cex = 0.5, col=4)
3548 legend(0, 9.5, legend = c('Sepal.length', 'Sepal.width',
3549           'Petal.length', 'Petal.width'),
3550           col = 1:4, lty = 1, lwd = 2, bty = 'n',
3551           y.intersp = 0.8, cex = 1.2)
3552 text(25, -1, 'Setosa\u1-50', cex = 1.3)
3553 text(75, -1, 'Versicolor\u51-100', cex = 1.3)
3554 text(125, -1, 'Virginica\u101-\u150', cex = 1.3)
3555 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\u1-50", size = 15)
plt.text(57,-1, "Versicolor\u51-100", size = 15)
plt.text(107,-1, "Virginica\u101-150", size = 15)

plt.show()

```

3556     The first 50 are setosa, which are the smallest flower, and the last 50 for virginica  
 3557     that have the largest petals. Given the data of sepal and petal sizes, one may  
 3558     correctly detect their corresponding flower species. However, not all the flowers have

3560 the same size. This makes the species detection more difficult. The random forest  
 3561 algorithm can make the detection with a small error. We make an RF experiment  
 3562 with this dataset. We randomly select 120 rows of data as the training data to  
 3563 obtained a trained RF model with 500 decision trees, use the remained 30 rows  
 3564 as the new data for our detection. These decision trees to vote on each of the 30  
 3565 rows of the new data. The most species votes a row receives from the trained RF  
 3566 trees is the RF detected result. By default, the decision trees in an RF algorithm  
 3567 are built randomly, the RF result from each run based on the same data may be  
 3568 slightly different. The following computer code shows the RF run and its result of  
 3569 this experiment for Fisher's iris data. The code also generates Fig. 3.9.

```

3570 #R code: RF prediction using the Fisher iris data
3571 #install.packages("randomForest")
3572 library(randomForest)
3573 #set.seed(8) # run this line to get the same result
3574 #randomly select 120 observations as training data
3575 train_id = sort(sample(1:150, 120, replace = FALSE))
3576 train_data = iris[train_id, ]
3577 dim(train_data)
3578 #[1] 120 5
3579 #use the remaining 30 as the new data for prediction
3580 new_data = iris[-train_id, ]
3581 dim(new_data)
3582 #[1] 30 5
3583
3584 #train the RF model
3585 classifyRF = randomForest(x = train_data[, 1:4],
3586                             y = train_data[, 5], ntree = 800)
3587 classifyRF #output RF training result
3588 #Type of random forest: classification
3589 #Number of trees: 800
3590 #No. of variables tried at each split: 2
3591
3592 #OOB estimate of error rate: 4.17%
3593 #Confusion matrix:
3594 #
3595 #          setosa versicolor virginica class.error
3596 #setosa      41          0          0  0.000000000
3597 #versicolor   0         34          2  0.055555556
3598 #virginica    0          3         40  0.06976744
3599
3600 plot(classifyRF) #plot the errors vs RF trees Fig. 9.9a
3601
3602 classifyRF$importance #classifyRF$ has many outputs
3603 #          MeanDecreaseGini
3604 #Sepal.Length      8.313131
3605 #Sepal.Width       1.507188
3606 #Petal.Length     31.075960
3607 #Petal.Width      38.169763
3608 varImpPlot(classifyRF) #plot the importance result Fig. 9.9b
3609
3610 #RF prediction for the new data based on the trained trees
3611 predict(classifyRF, newdata = new_data[,1:4])
3612 # 2           4           10          11          12          14
3613 #setosa      setosa      setosa      setosa      setosa      setosa

```

```

3613
3614 #It got two wrong: 120 versicolor and 135 versicolor
3615
3616 #Another version of the randomForest() command
3617 anotherRF = randomForest(Species ~ .,
3618                         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. 1. 2. 1. 1. 2. 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()
```

3620

```
#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()
```

3621

3622 The following explains the RF output data and figures.

- Confusion matrix: This matrix displays correct and wrong classifications based on the training data. The columns are truth and rows are RF classifications. The diagonals are corrected classifications, and off diagonals are the wrong ones. The 41 setosas are all correctly classified. Among the 37 versicolors, 34 are correctly classified, but 3 are wrongly classified as virginica. Among the 37 versicolors, 34 are correctly classified, but 3 are wrongly classified as virginica.
- Classification error: The last column in the confusion matrix is the classification error which is equal to the sum of the off-diagonal elements in that row divided by the sum of the entire row, i.e., the ratio of the number of wrong classifications to the total number of classifications. The errors shows how good is the trained RF model.
- OOB estimate of error rate: OOB stands for out-of-bag, a term for a statistical sampling process that puts a part of the data for training the RF model and the remainder for validation. This remainder is referred to as the OOB set

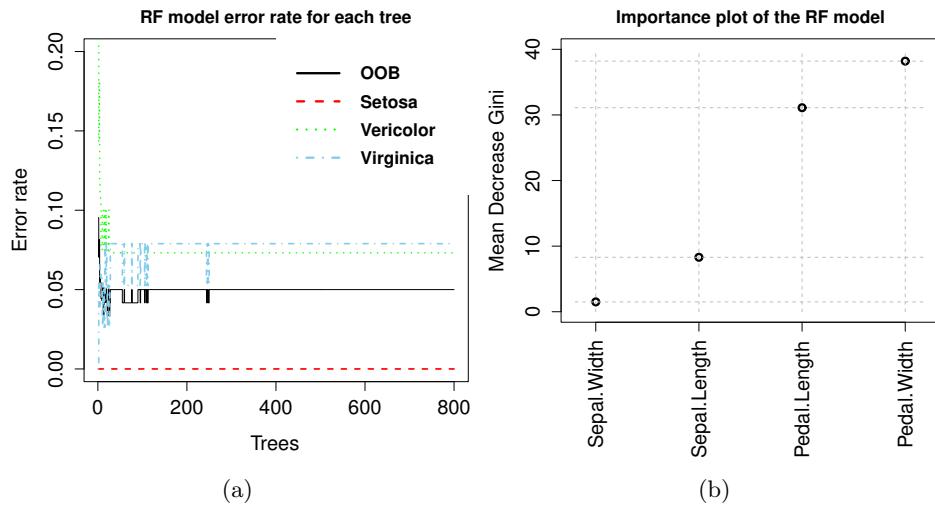


Fig. 3.9 (a) RF errors vs. trees. (b) The RF's importance plot.

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

3663       of the 9th day from today (e.g., sunny, rainy, or cloudy) for 30 different locations  
 3664       in a country.

### 3665       3.3.2 RF regression for the daily ozone data of New York City

---

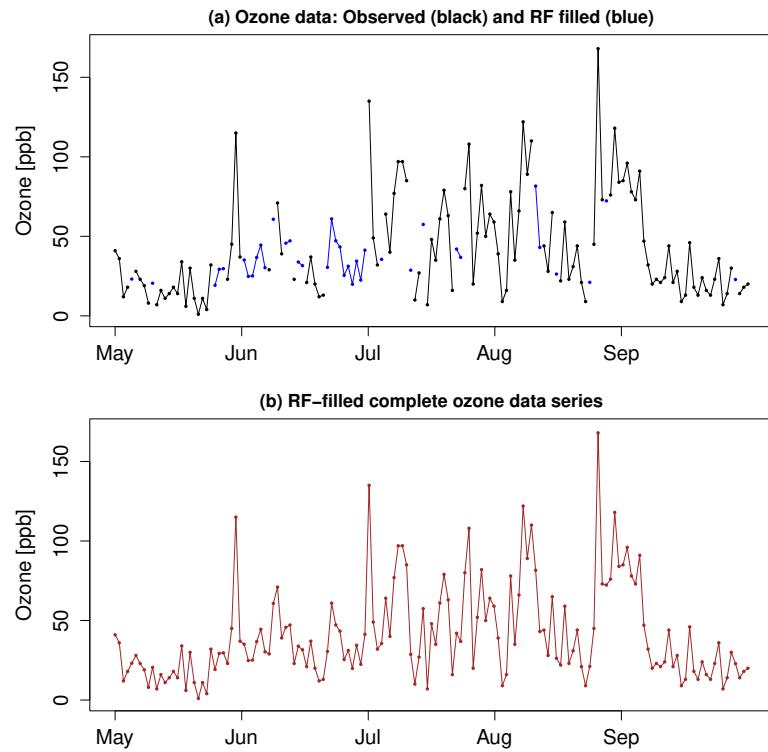
3667       Similar to Fisher's iris data, the daily air quality data, the file name `airquality`,  
 3668       of New York City (NYC) from 1 May 1973 - 30 September 1973 (153 days) is an-  
 3669       other benchmark data for the random forest algorithms. The dataset contains the  
 3670       following weather parameters: Ozone concentration in the ground-level air, cumu-  
 3671       lative solar radiation, average wind speed, and daily maximum air temperature.  
 3672       When the ozone concentration is higher than 86 parts per billion (ppb), the air is  
 3673       considered unhealthy. The ozone concentration is related to solar radiation, wind  
 3674       and temperature. The following list provides more information on this `airquality`  
 3675       dataset:

- 3676       • The ozone data in ppb observed between 1300 and 1500 hours at Roosevelt Island  
           3677       of the New York City. Among the 153 days, 37 had no data and are denoted by  
           3678       NA. The data range is 1 - 168 ppb. The maximum ozone level of this dataset  
           3679       is 168 ppb, occurred on 25 August 1973, when the cumulative solar radiation  
           3680       was high at 238 Lang, the average wind speed low at 3.4 mph, and Tmax is  
           3681       moderate 81°F.
- 3682       • Cumulative solar radiation from 0800 to 1200 hours with units in Langley (Lang  
           3683       or Ly) (1 Langley = 41,868 Watt·sec/m<sup>2</sup>, or 1 Watt/m<sup>2</sup> = 0.085985 Lang-  
           3684       ley/hour) in the lightwave length range 4,000 - 7,700 Angstroms at Central  
           3685       Park of the New York City. The data range is 7 - 334 Lang.
- 3686       • Average wind speed in miles per hour (mph) at 0700 and 1000 hours at LaGuardia  
           3687       Airport, less than 10 km from Central Park. The data range is 1.7 - 20.7 mph.
- 3688       • Maximum daily temperature Tmax in °F at La Guardia Airport. The data  
           3689       range is 56 - 97 °F. The daily Tmax data can also be downloaded from the  
           3690       NOAA NCEI website by an online search using the words: NOAA Climate Data  
           3691       Online LaGuardia Airport.

3692       While RF analysis for this `airquality` dataset can be for many purposes, our  
 3693       example is to fill the 37 missing ozone data with RF regression results. RF use the  
 3694       above four parameters and their data to build decision trees. Figure 3.10(a) shows  
 3695       the original 116 observed data points in black dots and lines, and the RF regression  
 3696       estimates for the 37 missing data (the blue dots and lines). A single blue dot means  
 3697       only an isolated day of missing data. A blue line indicates missing data in successive  
 3698       days. For better visualization, Fig. 3.10(b) shows the complete NYC daily ozone  
 3699       time series that join the observed data with the result data from an RF regression.

3700       Figure 3.10 may be generated by the following computer code.

```
3701 #R plot Fig. 9.10: RF regression for ozone data
3702 library(randomForest)
3703 airquality[1:2,] #use R's RF benchmark data "airquality"
3704 # 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.

```

3705 #1      41      190    7.4     67      5     1
3706 #2      36      118    8.0     72      5     2
3707 dim(airquality)
3708 #[1] 153   6
3709 ozoneRFreg = randomForest(Ozone ~ ., data = airquality,
3710                               mtry = 2, ntree = 500, importance = TRUE,
3711                               na.action = na.roughfix)
3712 #na.roughfix allows NA to be replaced by medians
3713 #to begin with when training the RF trees
3714 ozonePred = ozoneRFreg$predicted #RF regression result
3715 t0 = 1:153
3716 n1 = which(airquality$Ozone > 0) #positions of data
3717 n0 = t0[-n1] #positions of missing data
3718 ozone_complete = ozone_filled= airquality$Ozone
3719 ozone_complete[n0] = ozonePred[n0] #filled by RF
3720 ozone_filled = ozonePred #contains the RF reg result
3721 ozone_filled[n1] <- NA #replace the n1 positions by NA

```

```
3722 t1 = seq(5, 10, len = 153) #determine the time May - Sept
3723
3724 par(mfrow = c(2, 1))
3725 par(mar = c(3, 4.5, 2, 0.1))
3726 plot(t1, airquality$Ozone,
3727   type = 'o', pch = 16, cex = 0.5, ylim = c(0, 170),
3728   xlab = '', ylab = 'Ozone[ppb]', xaxt="n",
3729   main = '(a) Ozone data: Observed(black) and RF filled(blue)',
3730   col = 1, cex.lab = 1.3, cex.axis = 1.3)
3731 MaySept = c("May", "Jun", "Jul", "Aug", "Sep")
3732 axis(side=1, at=5:9, labels = MaySept, cex.axis = 1.3)
3733 points(t1, ozone_filled, col = 'blue',
3734   type = 'o', pch = 16, cex = 0.5)#RF filled data
3735
3736 #Plot the complete data
3737 par(mar = c(3, 4.5, 2, 0.1))
3738 plot(t1, ozone_complete,
3739   type = 'o', pch = 16, cex = 0.5, ylim = c(0, 170),
3740   xlab = '', ylab = 'Ozone[ppb]',
3741   xaxt="n", col = 'brown',
3742   main = '(b) RF-filled complete ozone data series',
3743   cex.lab = 1.3, cex.axis = 1.3)
3744 MaySept = c("May", "Jun", "Jul", "Aug", "Sep")
3745 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()
```

3747

```
#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()
```

3748

3749  
3750

### 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

3758 rows (31%) virginica. The first branch of the decision tree grows from a condition

$$\text{Petal Length} < 2.5 \text{ [cm]}. \quad (3.44)$$

3759 If “yes”, this iris is setosa. All the 41 setosa flowers (34% of the entire training  
 3760 data) have been detected. This decision is clearly supported by the real petal length  
 3761 data (the green line) shown in Fig. 3.8, which shows that only setosa has petal length  
 3762 less than 2.5.

3763 If “no”, then two possibilities exist. This allows us to grow a new branch of the  
 3764 tree. The condition for this branch is

$$\text{Petal Width} < 1.8 \text{ [cm]}. \quad (3.45)$$

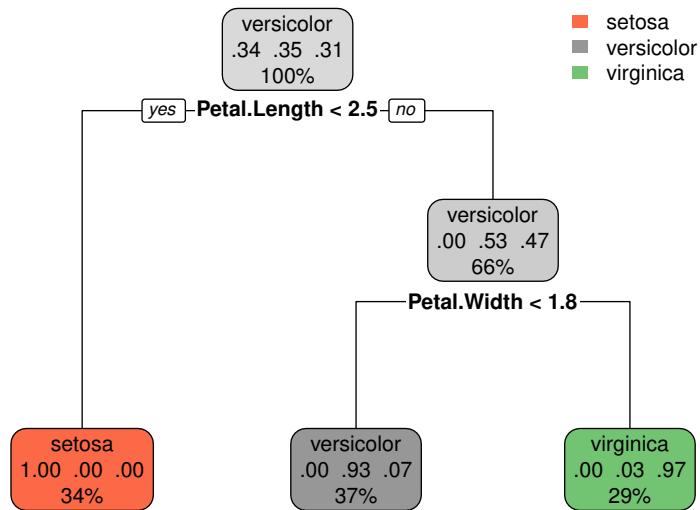
3765 This condition determines whether a flower is versicolor or virginica. After the  
 3766 isolation of 41 setosa flowers, the remaining 79 flowers are 42 versicolor (53%)  
 3767 and 37 virginica (47%). The “yes” result of condition (3.45) leads to versicolor.  
 3768 This result is 93% correct and 7% incorrect. The “no” result implies virginica. This  
 3769 conclusion is 97% correct, and 3% incorrect. These conclusions are supported by the  
 3770 real petal width data (the blue line) shown in Fig. 3.8. The petal width of versicolor  
 3771 iris flowers is in general less than 1.8 [cm]. However, the data have some fluctuations  
 3772 in both the versicolor and virginica sections in Fig. 3.8. These fluctuations lead to  
 3773 the errors of decision.

3774 The entire RF process for our iris species example had grown 800 such decision  
 3775 trees using the 120 rows of training data. Different trees use different branch  
 3776 conditions. RF algorithms grow these trees following various kinds of optimization  
 3777 principles, which involve some tedious mathematics not covered here.

3778 Figure 3.11 may be generated by the following computer code.

```
3779 #R plot of Fig. 9.11: A tree in a random forest
3780 #install.packages('rpart')
3781 library(rpart)
3782 #install.packages('rpart.plot')
3783 library(rpart.plot)
3784
3785 setwd('/~/climstats')
3786 setEPS() #Plot the data of 150 observations
3787 postscript("fig0911.eps", width=6, height=6)
3788 par(mar = c(0, 2, 1, 1))
3789 iris_tree = rpart( Species ~ . , data = train_data)
3790 rpart.plot(iris_tree,
3791   main = 'An decision tree for the RF training data')
3792 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()
```

3794

## 3.4 Neural network and deep learning

3795

3796 A neural network (NN) consists of a series of data fitting based on the training data  
3797 and an application of the fitted NN model for test data. The NN outputs categorical  
3798 predictions, such as a sunny or rainy day, or numerical predictions as a regression.  
3799 NN is also known as an artificial neural network (ANN) or simulated neural network  
3800 (SNN). It is a popular machine learning method, and is a fundamental building  
3801 block the deep learning algorithms that often involve the data fitting of multiple  
3802 layers, referred to as hidden layers.

3803 Why is it called neural network? How does a data fitting process have anything  
3804 to do with “neural” and/or “network”? Artificial neurons were first proposed by  
3805 Warren Sturgis McCulloch (1898 - 1969), an American neurophysiologist, and Wal-  
3806 ter Harry Pitts (1923 - 1969), an American logician, in their 1943 paper entitled “  
3807 A logical calculus of ideas immanent in nervous activity.” This mathematical paper  
3808 used terms “neuron”, “action”, “logic expression”, and “net”, which are among  
3809 the keywords in the modern NN writings. The ten theorems of the paper formu-  
3810 late a suite of logic expressions based on data. The word “neuron” was more a  
3811 graphic indication for actions and logic expressions than a biological reference. A  
3812 layman’s understanding of an NN machine learning is often put incorrectly toward  
3813 biological neurons and a biological neural network. Therefore, ANN may be a more  
3814 appropriate term for NN to avoid confusion.

3815 This book attempts to make such a brief NN introduction that you can use  
3816 our R or Python code for your data and objectives, interpret your NN computing  
3817 results, and understand the basic principles of mathematical formulations of an NN  
3818 algorithm. However, we do not attempt to derive mathematical details of the NN  
3819 theory. This section has two sub-sections: (i) A simple NN example of a decision  
3820 system, and (ii) An example of NN prediction using the benchmark data of Fisher’s  
3821 iris flowers.

3822  
3823

### 3.4.1 An NN model for an automated decision system

---

3824  
3825

#### 3.4.1.1 An overall idea of a neural network decision system for recruitment

3826 A senior human resource manager of an IT company recruited three new employ-  
3827 ees among six candidates. Her hiring decisions were made based on the technical  
3828 knowledge scores (TCS) and communication skills scores (CSS) given in Table 9.1.  
3829 TCS and CSS were the interview results from company’s working groups. A ju-  
3830 nior human resource manager wishes to use an NN model to help him make his  
3831 recruitment decisions consistent with those made by the senior manager.

3832 The junior manager has received the following scores for the three new job appli-  
3833 cants A, B, and C, whose TKS and CSS scores are as follows: TKS (30, 51, 72), and  
3834 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

3835 recruited when either TKS score or CSS score is high. Thus, this junior manager  
 3836 has an easy decision for Candidate A whose CSS 85 is high, so to be recruited.  
 3837 However, the decision for Candidate B is difficult. Neither TKS 51 nor CSS 51 of  
 3838 Candidate B is high, but both are higher than those of the rejected, and further the  
 3839 sum of TKS and CSS is 102, lower than the minimum of the totals of the recruited,  
 3840 110. Should the junior manager hire Candidate B? The decision for Candidate C is  
 3841 also difficult. Candidate C seems weaker than the recruited case TKS 30 and CSS  
 3842 80, but not too much weaker. Should the junior manager hire Candidate C? The  
 3843 following NN computer code may learn from the data of the senior manager and  
 3844 suggest an NN decision that may serve as a useful reference for the junior manager.

#### 3.4.1.2 An NN code, results, and their interpretation

```

3845 #R code for NN recruitment decision and Fig. 9.12
3846 #Ref: https://www.datacamp.com/tutorial/neural-network-models-r
3847 TKS = c(20,10,30,20,80,30)
3848 CSS = c(90,20,40,50,50,80)
3849 Recruited = c(1,0,0,0,1,1)
3850 #make a data frame for the NN function neuralnet
3851 df = data.frame(TKS, CSS, Recruited)
3852
3853 require(neuralnet) # load 'neuralnet' library
3854 # fit neural network
3855 set.seed(123)
3856 nn = neuralnet(Recruited ~ TKS + CSS, data = df,
3857                 hidden = 5, act.fct = "logistic",
3858                 linear.output = FALSE)
3859 plot(nn) #Plot Fig 9.12: A neural network
3860
3861 TKS=c(30,51,72) #new data for decision
3862 CSS=c(85,51,30) #new data for decision
3863 test=data.frame(TKS,CSS)
3864 Predict=neuralnet::compute(nn,test)
3865 Predict$net.result #the result is probability
3866 #[1,] 0.99014936
3867 #[2,] 0.58160633
3868 #[3,] 0.01309036
3869
3870 # Converting probabilities into decisions
3871 ifelse(Predict$net.result > 0.5, 1, 0) #threshold = 0.5
3872 #[1,] 1
3873 #[2,] 1
3874 #[3,] 0
3875
3876 #print bias and weights
3877 nn$weights[[1]][[1]]
3878

```

```
3879 #print the last bias and weights before decision
3880 nn$weights[[1]][[2]]
3881 #print the random start bias and weights
3882 nn$startweights
3883 #print error and other technical indices of the nn run
3884 nn$result.matrix #results data
```

```

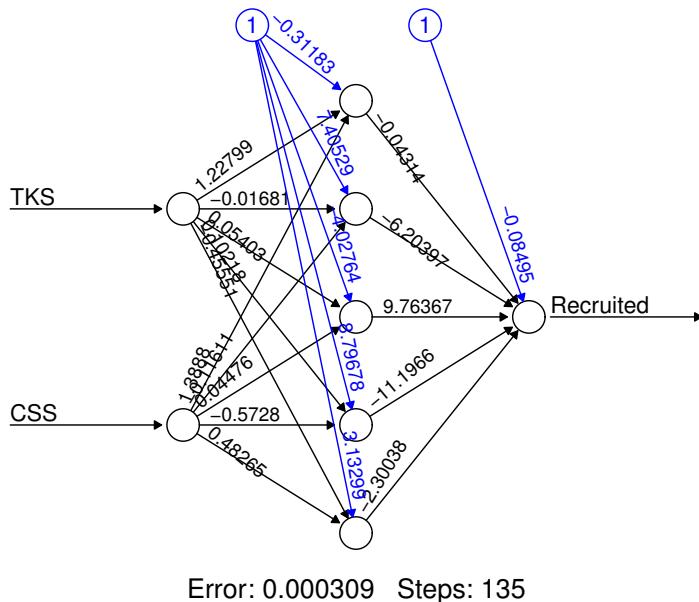
#Python code for the NN recruitment decision and Fig. 9.12
#imports to build neural network model and visualize
#You need to install the following Python environments:
#tensorflow and keras-models
#!pip3 install ann_visualizer
#!pip3 install keras-models
#from a terminal: pip install tensorflow
from 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\u2225Decision")
#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:

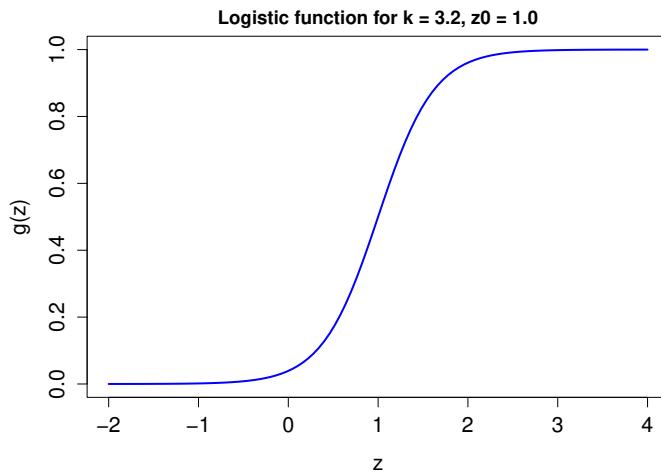
$$z_j = \sum_{i=1}^n w_{ij}x_i + b_j, \quad (3.46)$$

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

$$g(z) = \frac{1}{1 + \exp(-k(z - z_0))}, \quad (3.47)$$

where  $k$  is called the logistic growth rate, and  $z_0$  is the midpoint. This function can also be written as

$$g(z) = \frac{1}{2} + \frac{1}{2} \tanh\left(-\frac{k(z - z_0)}{2}\right), \quad (3.48)$$



**Fig. 3.13** Logistic function with growth rate  $k = 3.2$  and midpoint  $y_0 = 1.0$ .

3898 where the tanh function is defined as

$$\tanh(t) = \frac{\exp(t) - \exp(-t)}{\exp(t) + \exp(-t)}. \quad (3.49)$$

3899 A curve for a logistic function is shown in Fig. 3.13. The curve has a property that  
 3900  $f(-\infty) = 0$  and  $f(\infty) = 1$ . This property makes the logistic activation function  
 3901 useful for categorical assignment: False for 0 and True for 1.

3902 Figure 3.13 can be plotted by the following computer code.

```
3903 #R plot Fig. 9.13: Curve of a logistic function
3904 y = seq(-2, 4, len = 101)
3905 k = 3.2
3906 y0 = 1
3907 setEPS() #Automatically saves the .eps file
3908 postscript("fig0913.eps", height=5, width=7)
3909 par(mar = c(4.2, 4.2, 2, 0.5))
3910 plot(y, 1/(1 + exp(-k*(y - y0))),
3911       type = 'l', col = 'blue', lwd = 2,
3912       xlab = 'y', ylab = 'f(y)',
3913       main = 'Logistic function for k = 3.2, y0 = 1.0',
3914       cex.lab = 1.3, cex.axis = 1.3)
3915 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()

```

3916

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. (3.46) 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

3947 than six groups of training data. Specifically, we will use 75 groups of the Fisher  
 3948 iris flower data.

### 3.4.1.3 Dissect NN results and errors

3950 Logistic regression fits weighted data to a logistic function. Logistic regression can  
 3951 be a standalone statistics course or a chapter in a statistics book. Interested readers  
 3952 are referred to those books and courses.

3953 In addition, there are other types of activation functions, such as linear, rectified  
 3954 linear, leaky rectified linear, exponential linear, and SoftPlus. You can find these  
 3955 functions from the specialized ML books or online.

3956 Usually a software package of neural network can output the weights and bias,  
 3957 such as using the following R command

3958 `nn$weights[[1]][[1]]}`

3959

3960 in the R NN package `neuralnet`. R can also output the initial assignment of weights  
 3961 and bias

3962 `nn$weights[[1]][[2]]}`

3963

3964 The initial weights and biases are randomly assigned, which makes the NN result  
 3965 different for different runs when having insufficient amount of data, or insufficient  
 3966 number of hidden layers.

3967 An NN computer package can also output many other modeling results, such as  
 3968 an R command

3969 `nn$result.matrix`

3970

3971 These output data can be used to analyze the quality of your trained model. So, it  
 3972 is unfair to say that NN is a blackbox. Yet, analyzing a trained model for a complex  
 3973 NN is very difficult and requires some nontrivial mathematical preparation.

3974 Most NN users may not have the mathematical ability to analyze the trained  
 3975 model based on mathematical theories. Instead, they apply the model to the new  
 3976 data and see if the NN prediction makes sense based on their experience or common  
 3977 sense. If not, train the model again with different parameters, or even feed more  
 3978 training data.

### 3.4.2 An NN prediction of iris species

---

3979  
 3980 We wish to use the Fisher iris data to show an example of NN with a sufficiently  
 3981 large size of training data and a highly reliable result. The Fisher iris data is a  
 3982 150 × 5 matrix with three iris species: 50 setosa, 50 virginica, and 50 versicolor, as  
 3983 described in the random forest section (Fisher 1936). We wish to use a percentage of  
 3984

3985 the data as the training data and the rest as the test data. The following computer  
 3986 code is for 50% of the data to be the training data, i.e.,  $p = 0.5$  in the code. We  
 3987 choose 10 neurons, and 10 hidden layers. The NN prediction result shows that NN  
 3988 correctly predicted all the 27 setosa irises, got 19 correct among the 20 versicolor  
 3989 irises but got one versicolor incorrectly as virginica, and got 26 correct among the  
 3990 28 virginica irises but got 2 incorrectly as versicolor. Since the 75 groups of training  
 3991 data are plenty, compared to the 6 groups of training data for the hiring decision,  
 3992 the trained NN model for the iris flowers does not vary much for different runs.  
 3993 The prediction results are quite reliable. Using more neurons and hidden layers may  
 3994 also help achieve reliable results. Although the results from different runs of the NN  
 3995 code do not change much, small differences still exist. For example, you may get a  
 3996 result of predicting 27 virginica.

3997 When you change  $p = 0.8$ , you will have 120 rows of training data and 30 rows of  
 3998 test data. This is the same as what we used in the random forest example section.  
 3999 This increase of training data size helps further stabilize the prediction results, i.e.,  
 4000 the results from different runs have little differences. You can make some runs of  
 4001 the code on your own computer, and see how many incorrect predictions are there.  
 4002 Our example runs shown below show only one or two incorrect predictions among  
 4003 the 30 rows of test data. For example, in one run, the 9 setosa and 12 versicolor  
 4004 irises were correctly predicted, and the 9 virginica irises had 8 correct predictions  
 4005 and one wrong. The level of the NN prediction accuracy seems similar to that of  
 4006 the random forest prediction for this particular dataset.

```
4007 #R NN code for the Fisher iris flower data
4008 #Ref: https://rpubs.com/vitorhs/iris
4009 data(iris) #150-by-5 iris data
4010 #attach True or False columns to iris data
4011 iris$setosa = iris$Species == "setosa"
4012 iris$virginica = iris$Species == "virginica"
4013 iris$versicolor = iris$Species == "versicolor"
4014 p = 0.5 # assign 50% of data for training
4015 train.idx = sample(x = nrow(iris), size = p*nrow(iris))
4016 train = iris[train.idx,] #determine the training data
4017 test = iris[-train.idx,] #determine the test data
4018 dim(train) #check the train data dimension
4019 #[1] 75 8
4020
4021 #training a neural network
4022 library(neuralnet)
4023 #use the length, width, True and False data for training
4024 iris.nn = neuralnet(setosa + versicolor + virginica ~
4025           Sepal.Length + Sepal.Width +
4026           Petal.Length + Petal.Width,
4027           data = train, hidden=c(10, 10),
4028           rep = 5, err.fct = "ce",
4029           linear.output = F, lifesign = "minimal",
4030           stepmax = 1000000, threshold = 0.001)
4031
4032 plot(iris.nn, rep="best") #plot the neural network
4033
4034 #Prediction for the rest data
```

```
4035 prediction = neuralnet::compute(iris.nn, test[,1:4])
4036 #prediction$net.result is 75-by-3 matrix
4037 prediction$net.result[1:2,]
4038 # [,1]          [,2]          [,3]
4039 #2    1 3.531638e-09 6.745692e-137
4040 #4    1 3.527872e-09 6.756521e-137
4041
4042 #find which column is for the max of each row
4043 pred.idx <- apply(prediction$net.result, 1, which.max)
4044 pred.idx
4045 #2 4 6 8 9 10
4046 #1 1 1 1 1 1
4047
4048 #Assign 1 for setosa, 2 for versicolor, 3 for virginica
4049 predicted <- c('setosa', 'versicolor', 'virginica')[pred.idx]
4050 predicted[1:6] #The prediction result
4051 #[1] "setosa" "setosa" "setosa" "setosa" "setosa"
4052
4053 #Create confusion matrix: table(prediction,observation)
4054 table(predicted, test$Species)
4055 #predicted      setosa versicolor virginica
4056 #setosa          27      0      0
4057 #versicolor      0      19      2
4058 #virginica       0      1      26
```



4060

### 3.5 Chapter summary

4061

4062 This chapter has described four methods of machine learning: K-means, support  
4063 vector machine (SVM), random forest (RF), and neural network (NN). In the de-  
4064 scription, we have used the following datasets: the daily weather data at the Miami  
4065 International Airport, the United States, the daily air quality data of New York  
4066 City, and R. A. Fisher's iris flower data of species. We have provided both R and  
4067 Python codes for these algorithms and their applications examples. The following  
4068 provides a brief summary of our descriptions about the four methods.

- 4069 (i) K-means clustering: Simply speaking, this clustering method can fairly divide  
4070  $N$  identical candies on a table for  $K$  kids according to the candy locations.  
4071 It is normally used as an unsupervised learning method. The K-means algo-  
4072 rithm minimizes the total within cluster sum of squares (tWCSS) through  
4073 iterations. As an example, we applied the K-means clustering method to  
4074 the data of the daily minimum surface air temperature  $T_{min}$  [in  $^{\circ}\text{C}$ ] and  
4075 the direction of the fastest 2-minute wind in a day  $WDF2$  [in degrees] of the  
4076 Miami International Airport. Our tests suggested that the data have two  
4077 clusters, one corresponding to the prevailing wind from east, and another  
4078 from west.
- 4079 (ii) Support vector machine for the maximum separation of different sets: The  
4080 SVM algorithm is built on the principle of the maximum differences among  
4081 the labelled groups of data. The maximum separation of different groups is  
4082 measured by the distance between positive and negative hyperplanes. SVM  
4083 can also predict the labels of the new unlabeled data. SVM is usually used  
4084 as a tool of supervised learning.
- 4085 (iii) Random forest of decision trees: An RF model is a set of decision trees which  
4086 form a "forest of decisions." It is usually used as a tool of supervised learn-  
4087 ing. It has two steps: training and prediction. RF can make both classifi-  
4088 cation and regression. A benchmark RF example is used in our RF classi-  
4089 fication description: The separation and prediction of the R.A. Fisher iris  
4090 species dataset. We have also included an example on RF regression using  
4091 the daily air quality data of New York City.
- 4092 (iv) Neural network of multiple data fitting: An NN model is a series of data  
4093 fitting for a given activation function. NN can also make classifications and  
4094 numerical predictions as a regression. The logistic function is a popular  
4095 activation function, a smooth transition between category 0 to category 1.  
4096 Because of randomness is involved in the NN algorithm, different runs of R  
4097 or Python code may yield different results. We have included an example  
4098 of using the NN model to help make a hiring decision using the existing  
4099 data. To compare with the RF model, we have also applied the NN model  
4100 to the R.A. Fisher flower dataset of iris species.

4101

## References and Further Readings

- 4102 [1] Boehmke, B., and B. Greenwell, 2019: *Hands-on Machine Learning with R*.  
4103 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/>.

- 4104  
4105 [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.

- 4107  
4108 [3] Fisher, R.A., 1936: The use of multiple measurements in taxonomic problems.  
4109 *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 (1890-1962), 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/](https://www.angela1c.com/projects/iris_project/the-iris-dataset/)

- 4110  
4111 [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.

- 4114  
4115 [5] McCulloch, W. and W. Pitts, 1943: A logical calculus of ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics* 5, 115 -133.

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.

4117

- 4118 [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,  
4119 USA.  
4120

This is a very famous machine learning book, which has in depth descriptions of the commonly used ML algorithms.

4121

4122

## Exercises

4123

- 4124 **3.1** Use the K-means method to make a cluster analysis for the following five  
4125 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  
4126  $K = 2$ . Plot a figure similar to Fig. 9.1. What is the final tWCSS equal to?  
4127  
4128 **3.2** Given the following three points  $P_1(1, 1), P_2(2, 2), P_3(2, 3)$  and given  $K =$   
4129 2, make a K-means cluster analysis following the method presented in Sub-  
Section 9.1.1. Plot your K-means clustering result.  
4130  
4131 **3.3** Use the K-means method to make a cluster analysis for the daily TMIN and  
4132 WDF2 data of the Miami International Airport in 2015, following the method  
4133 presented in Sub-sections 9.1.3. Here, TMIN is the daily minimum temper-  
4134 ature, and WDF2 denotes the direction [in degrees] of the fastest 2-minute  
4135 wind in a day. You can obtain the data from the NOAA Climate Data Online  
website, or from the file

4136 `MiamiIntlAirport2001_2020.csv`

4137 in the `data.zip` file downloadable from the book website

4138 [www.climatestatistics.org](http://www.climatestatistics.org)

4139

- 4140 **3.4** Use the K-means method to make a cluster analysis for the daily TMAX and  
4141 WDF2 data of the Miami International Airport in 2015.  
4142 **3.5** Use the K-means method to make a cluster analysis for the daily TMIN and  
4143 Tmax data of the Miami International Airport in 2015.  
4144 **3.6** Use the K-means method to make a cluster analysis for the daily TMIN and  
4145 PRCP data of the Miami International Airport in 2015. Here, PRCP denotes  
4146 the daily total precipitation.

- 4147   **3.7** Identify two climate parameters of your interest, find the corresponding data  
4148   online, and make a K-means cluster analysis similar to the previous problem  
4149   for your data.
- 4150   **3.8** Following the method presented in Sub-section 9.2.1., make an SVM analysis  
4151   for the following five points in a 2D space:  $P_1(1, 1)$ ,  $P_2(2, 2)$ ,  $P_3(2, 3)$ ,  $P_4(3, 4)$ ,  $P_5(4, 4)$ .  
4152   The first three points are labeled 1 and the last two are labeled 2. What are  
4153    $w$ ,  $b$  and  $D_m$ ? What points are the support vectors?
- 4154   **3.9** Two new points are introduced to the previous problem:  $Q_1(1.5, 1)$  and  $Q_2(3, 3)$ .  
4155   Use the SVM trained in the previous problem to find out which point belongs  
4156   to what category.
- 4157   **3.10** From the Internet, download the historical daily data of minimum tempera-  
4158   ture (TMIN) and average wind speed (AWND) for a month and a location of  
4159   your interest. Label your data as 1 if the daily precipitation is greater than  
4160   0.5 millimeter, and 0 otherwise. Make an SVM analysis for your labeled data.
- 4161   **3.11** SVM forecast: From the Internet, download the historical daily data of min-  
4162   imum temperature and sea level pressure for a month and a location of your  
4163   interest. Label your data as 1 if the total precipitation of the next day is  
4164   greater than 0.5 millimeter, and 0 otherwise. Make an SVM analysis for your  
4165   labeled data. Given the daily data of minimum temperature and sea level pres-  
4166   sure of a day in the same month but a different year, use your trained SVM  
4167   model to forecast whether the next day was rainy or not, i.e., to determine  
4168   whether the next day is 1 or 0.
- 4169   **3.12** In order to improve the accuracy of your prediction, can you use the data of  
4170   multiple years to train your SVM model? Perform numerical experiments and  
4171   show your results.
- 4172   **3.13** The first 50 in the 150 rows of the R.A. Fisher iris data are for setosa, 51-100  
4173   are for versicolor, and 101-150 are for virginica, use the data 1-40, 51-90, and  
4174   101-140 to train an RF model, following the method in Sub-Section 9.3.1.  
4175   Use the RF model to predict the species of the remaining data of lengths  
4176   and widths of petal and sepal. You can download the R.A. Fisher dataset  
4177   `iris.csv` from the book website or use the data already built in R or Python  
4178   software packages.
- 4179   **3.14** For the 150 rows R.A. Fisher iris data, use only 20% of the data from each  
4180   species to train your RF model. Select another 10% of the riris data of lengths  
4181   and widthes as the new data for prediction. Then use the RF model to predict  
4182   the species of the new data. Discuss the errors of your RF model and your  
4183   prediction.
- 4184   **3.15** Plot a decision tree like Fig. 3.11 for the previous exercise problem.
- 4185   **3.16** For the same R.A. Fisher dataset, design your own RF training and predic-  
4186   tion. Discuss the RF prediction accuracy for this problem.
- 4187   **3.17** RF forecast: From the Internet, download the historical daily data of mini-  
4188   mum temperature and sea level pressure for a month and a location of your  
4189   interest. Label your data as 1 if the total precipitation of the next day is  
4190   greater than 0.5 millimeter, and 0 otherwise. Make an RF analysis for your

- 4191 labeled data. Given the daily data of minimum temperature and sea level  
4192 pressure of a day in the same month but a different year, use your trained  
4193 RF model to forecast whether the next day was rainy or not. Is your forecast  
4194 accurate? How can you improve your forecast?
- 4195 **3.18** Find a climate data time series with missing data and fill in the missing data  
4196 using the RF regression method described in Sub-Section 9.3.2.
- 4197 **3.19** The first 50 in the 150 rows of the R.A. Fisher iris data are for setosa, 51-100  
4198 are for versicolor, and 101-150 are for virginica, use the data 1-40, 51-90, and  
4199 101-140 to train an NN model, following the method in Sub-Section 9.4.2.  
4200 Use the NN model to predict the species of the remaining data of lengths and  
4201 widths of petal and sepal.
- 4202 **3.20** For the 150 rows R.A. Fisher iris data, use only 20% of the data from each  
4203 species to train your NN model. Select another 10% of the riris data of lengths  
4204 and widthes as the new data for prediction. Then use the NN model to predict  
4205 the species of the new data. Discuss the errors of your NN model and your  
4206 prediction.
- 4207 **3.21** NN forecast: From the Internet, download the historical daily data of mini-  
4208 mum temperature and sea level pressure for a month and a location of your  
4209 interest. Label your data as 1 if the total precipitation of the next day is  
4210 greater than 0.5 millimeter, and 0 otherwise. Make an NN analysis for your  
4211 labeled data. Given the daily data of minimum temperature and sea level  
4212 pressure of a day in the same month but a different year, use your trained  
4213 NN model to forecast whether the next day was rainy or not. Is your forecast  
4214 accurate? Compare your NN forecast with your RF and SVM forecasts.
- 4215 **3.22** Design a machine learning project for yourself or others. What are your train-  
4216 ing data? What are your test data? What is your training model? What is  
4217 your training model error? How would you assess your prediction error?

## Matrix Applications to Regression Models

4220

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.

4237

### 4.1 Simple linear regression

4238

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.

4245

4246

#### 4.1.1 Temperature lapse rate and an approximately linear model

4247

4248

4249

4250

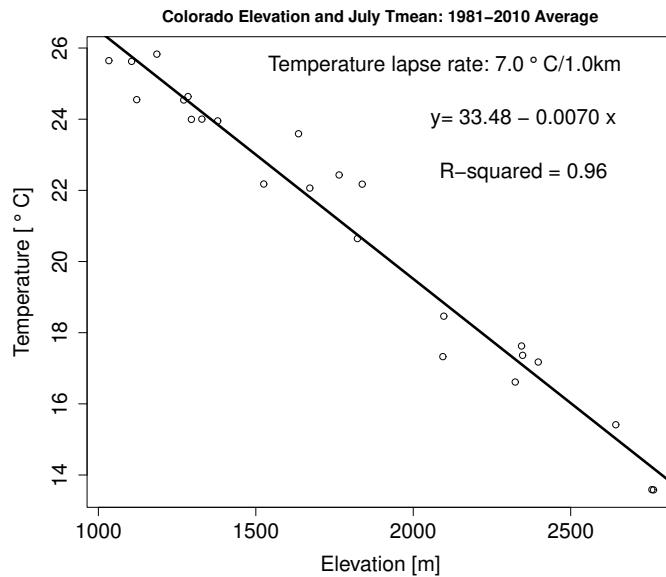
It is known that as an approximation, in a mountain terrain temperature decreases linearly according to elevation. The decrease rate is called the temperature *lapse rate* in meteorology. The linear relationship and the lapse rate are shown in Fig. 4.1 for Colorado, which is a Rocky Mountains state of the United States and has

4251 a large elevation range from 1010 to 4401 meters, according to Colorado Tourism  
 4252 [www.colorado.com](http://www.colorado.com). The dots in the figure are the scatter plot of the data of elevation  
 4253 and 1981–2010 average July surface air temperature Tmean of the 24 Colorado  
 4254 stations in the United States Historical Climatological Network (USHCN) (Menne  
 4255 et al. 2009). The straight line of the figure corresponds to the linear model

$$y = 33.48 - 0.0070x. \quad (4.2)$$

4256 This equation is a result from regression of data, referred to as observations in  
 4257 climate science, or sample data in statistics.

4258 Thus, Figure 4.1 shows an exhibit of a scattered and complex data pattern re-  
 4259 gressed to a simpler form, in this case, a straight line model. This is called the simple  
 4260 linear regression process. The regression result has climate science interpretations,  
 4261 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.

4262 The horizontal axis of Fig. 4.1 is for elevation. The elevation data in units of  
 4263 meters for the 24 stations are as follows:

4264

4265 1671.5, 1635.6, 2097.0, 1295.4, 1822.7, 2396.9, 2763.0, 1284.7,  
 4266 1525.2, 1328.6, 1378.9, 2323.8, 2757.8, 1033.3, 1105.5, 1185.7,  
 4267 2343.9, 1764.5, 1271.0, 2347.3, 2094.0, 2643.2, 1837.9, 1121.7

4268 The vertical axis is for temperature, which is the 30-year average of the July daily

4269 mean temperature (Tmean) from 1981-2010, computed from the USHCN monthly  
 4270 data, which have been adjusted for the time of observation bias (TOB). Some  
 4271 stations had missing data denoted by -9999. When computing the 30-year average,  
 4272 the entries of -9999 were omitted. Thus, some averages were computed from fewer  
 4273 than 30 years. The resulting average temperature data in the unit of °C for the 24  
 4274 stations are as follows:

4275 22.064, 23.591, 18.464, 23.995, 20.645, 17.175, 13.582, 24.635,  
 4276 22.178, 24.002, 23.952, 16.613, 13.588, 25.645, 25.625, 25.828,  
 4277 17.626, 22.433, 24.539, 17.364, 17.327, 15.413, 22.174, 24.549

4278 With the above data, Figure 4.1 can be generated by the following computer code

```
4279 #R plot Fig. 4.1: Colorado temperature lapse rate
4280 x= c(
4281 1671.5, 1635.6, 2097.0, 1295.4, 1822.7, 2396.9, 2763.0, 1284.7,
4282 1525.2, 1328.6, 1378.9, 2323.8, 2757.8, 1033.3, 1105.5, 1185.7,
4283 2343.9, 1764.5, 1271.0, 2347.3, 2094.0, 2643.2, 1837.9, 1121.7)
4284 y= c(
4285 22.064, 23.591, 18.464, 23.995, 20.645, 17.175, 13.582, 24.635,
4286 22.178, 24.002, 23.952, 16.613, 13.588, 25.645, 25.625, 25.828,
4287 17.626, 22.433, 24.539, 17.364, 17.327, 15.413, 22.174, 24.549)
4288 setEPS() # save the .eps figure
4289 postscript("fig0401.eps", width = 8)
4290 par(mar=c(4.5,4.5,2.5,0.5))
4291 plot(x,y,
4292   xlab="Elevation [m]",
4293   ylab=expression("Temperature [~degree ~" "C]"),
4294   main="Colorado Elevation and July Tmean: 1981-2010 Average",
4295   cex.lab=1.5, cex.axis=1.5, cex.main = 1.2)
4296 reg=lm(y~x)
4297 reg
4298 #(Intercept)           x
4299 # 33.476216 -0.006982  #-7.0 degC/1km
4300 summary(reg)
4301 #R-squared:  0.9631
4302 abline(reg,lwd=3)
4303 text(2100, 25.5,
4304 expression("Temperature lapse rate: 7.0 ~degree ~" "C/1.0 km"),
4305 cex=1.5)
4306 text(2350, 24, "y= 33.48 - 0.0070 x", cex=1.5)
4307 text(2350, 22.5,"R-squared = 0.96", cex=1.5)
4308 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.\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"$y\u2022=\u202233.48\u2022-\u20220.0070^{\u2022}x$" % Rsqu,
        color= 'k', size = 20)
ax.text(2250, 23.5, r"$R\u2022-squared\u2022=\u2022%.2f$" % Rsqu,
        color= 'k', size = 20)
plt.show()

```

4309

4310     The Colorado July temperature lapse rate (TLR) 7.0 °C/1.0km is consistent with  
 4311     earlier studies using 104 stations on the west slope of Colorado for an elevation range  
 4312     of 1,500- 3,600 meters: 6.9 °C/1.0km (Fall 1997). For the annual mean temperature,  
 4313     the LTR is 6.0 °C/1.0km. These two results are comparable to an LTR study for the  
 4314     Southern Ecuadorian Andes in the elevation range of 2,610-4,200 meters (Cordova  
 4315     et al. 2016). The annual mean temperature TLR is 6.88 °C/1.0km.

4316     TLR may be applied to approximately predict the temperature of a mountain  
 4317     region at a given location, in case it is hard to maintain a weather station at some  
 4318     locations due to high elevation or complex terrain, while it is relatively easy to  
 4319     obtain the elevation data based on the Geographical Information System (GIS) or  
 4320     a digital elevation model (DEM) dataset.

---

### 4.1.2 Assumptions and formula derivations of the single variate linear regression

---

This subsection describes the statistical concepts and mathematical theory of the above linear regression procedure. We pay special attention to the assumptions about the linear regression model.

#### 4.1.2.1 Linear model and its assumptions

The linear model for TLR may be written as

$$Y = a + bx + \epsilon. \quad (4.3)$$

Here,

- (i)  $x$  is an independent variable, also called explanatory variable, and is deterministic (not random);
- (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  $x$ ;
- (iii)  $\epsilon$  is the random error term, which is assumed to have zero mean  $E[\epsilon] = 0$ , and constant variance  $\sigma^2$ :  $\text{Var}(\epsilon|x) = \sigma^2$ , and is uncorrelated with each other:  $\text{Cor}(\epsilon|x_1, \epsilon|x_2) = 0$  if  $x_1 \neq x_2$ ;
- (iv)  $a, b$  and  $\sigma$  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
- (v) The expected value of  $Y$  for a given  $x$  is  $a + bx$ , i.e.,

$$E[Y|x] = a + bx. \quad (4.4)$$

The linear model Eq. (4.3) at the given points  $x_i$  is

$$Y_i = a + bx_i + \epsilon_i, \quad i = 1, 2, \dots, n. \quad (4.5)$$

where both  $\epsilon_i$  and  $Y_i$  are random variables with

$$E[Y_i] = a + bx_i, \quad (4.6)$$

$$E[\epsilon_i] = 0, \quad (4.7)$$

$$\text{Var}(Y_i) = E[(Y_i - (a + bx_i))^2] = \sigma^2, \quad (4.8)$$

$$\text{Var}(\epsilon_i) = E[(\epsilon_i - 0)^2] = \sigma^2, \quad (4.9)$$

$$\text{Cor}(\epsilon_i, \epsilon_j) = \delta_{ij}, \quad (4.10)$$

where  $\delta_{ij}$  is the Kronecker delta

$$\delta_{ij} = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{if } i \neq j. \end{cases} \quad (4.11)$$

The last equation means that the error terms  $\epsilon_i$  at the observational points  $x_i$  are uncorrelated with each other.

4345 The observations at the given points  $x_i$  are  $y_i$ . Both are fixed values, such as tem-  
 4346 perature  $y_1 = 22.064^\circ\text{C}$  at the given elevation  $x_1 = 1671.5$  meters in the Colorado  
 4347 July temperature example. Here,  $y_1$  is a sample value for the random variable  $Y_1$ .  
 4348 The corresponding error datum is  $e_1$ , which is called a residual and is regarded a  
 4349 sample value of  $\epsilon_1$  with respect to the linear model.

#### 4350 4.1.2.2 Estimating $a$ and $b$ from data

4351 Given a set of sample data  $(x_i, y_i), i = 1, 2, \dots, n$ , the parameters  $a$  and  $b$  in the  
 4352 regression model can be estimated as  $\hat{a}, \hat{b}$ . Thus,  $\hat{a}, \hat{b}$  correspond to a particular  
 4353 dataset, i.e., another dataset  $(u_i, v_i), i = 1, 2, \dots, m$  will yield a pair of different es-  
 4354 timators that may be denoted by  $\hat{a}_{uv}, \hat{b}_{uv}$ . In contrast,  $a, b$  are the general notations  
 4355 of constant coefficients for a linear regression model.

4356 The estimated linear model may be written in the following way:

$$y_i = \hat{a} + \hat{b}x_i + e_i, \quad i = 1, 2, \dots, n, \quad (4.12)$$

4357 where

$$e_i = y_i - (\hat{a} + \hat{b}x_i) \quad (4.13)$$

4358 are called residuals of the linear model, and are equal to the observed values  $y_i$   
 4359 minus the predicted values

$$\hat{y}_i = \hat{a} + \hat{b}x_i. \quad (4.14)$$

4360 The mean of the above  $n$  equations (4.12) yields

$$\bar{y} = \hat{a} + \hat{b}\bar{x}, \quad (4.15)$$

4361 where

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}, \quad (4.16)$$

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{n}. \quad (4.17)$$

4362 Because of the unbiased model assumption  $E[\epsilon] = 0$ , we require that

$$\bar{e} = \frac{e_1 + e_2 + \dots + e_n}{n} = 0, \quad (4.18)$$

4363 equivalently,

$$\sum_{i=1}^n e_i = 0. \quad (4.19)$$

4364 Equations (4.15) and (4.18) give an estimate for  $a$  from data:

$$\hat{a} = \bar{y} - \hat{b}\bar{x}, \quad (4.20)$$

4365 but  $\hat{b}$  is to be estimated by

$$\hat{b} = \frac{\mathbf{y}_a \cdot \mathbf{x}_a}{\mathbf{x}_a \cdot \mathbf{x}_a}. \quad (4.21)$$

<sup>4366</sup> Here,

$$\mathbf{x}_a = \begin{bmatrix} x_{a,1} \\ x_{a,2} \\ \vdots \\ x_{a,n} \end{bmatrix} \quad \mathbf{y}_a = \begin{bmatrix} y_{a,1} \\ y_{a,2} \\ \vdots \\ y_{a,n} \end{bmatrix} \quad \mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}, \quad (4.22)$$

<sup>4367</sup> with

$$x_{a,i} = x_i - \bar{x}, \quad (4.23)$$

$$y_{a,i} = y_i - \bar{y}. \quad (4.24)$$

<sup>4368</sup> In the last two formulas,  $x_{a,i}$  and  $y_{a,i}$  may be regarded as anomaly data, hence  
<sup>4369</sup> with the subscript “a.” Climate scientists often uses the concept of anomaly data,  
<sup>4370</sup> departures from a climate normal.

<sup>4371</sup> The  $\hat{b}$  estimate formula (4.21) can be derived from the best fit condition which  
<sup>4372</sup> minimizes the mean square errors (MSE)

$$\text{MSE} = \frac{|\mathbf{e}|^2}{n} = \frac{\sum_{i=1}^n e_i^2}{n} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}, \quad (4.25)$$

<sup>4373</sup> equivalently minimizing the sum of the square errors (SSE)

$$\text{SSE} = |\mathbf{e}|^2 = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = n \times \text{MSE}, \quad (4.26)$$

<sup>4374</sup> where  $|\mathbf{e}|$  denotes the Euclidean length of vector  $\mathbf{e}$ . Thus, this method is also called  
<sup>4375</sup> the least square estimate. The least square condition is equivalent to the following  
<sup>4376</sup> orthogonality condition (See Fig. 4.2)<sup>1</sup>

$$\mathbf{e} \cdot \mathbf{x}_a = 0. \quad (4.27)$$

<sup>4377</sup> Inserting Eq. (4.20) into the linear model with data (4.12) leads to  
<sup>4378</sup>

$$y_i - \bar{y} = \hat{b}(x_i - \bar{x}) + e_i, \quad i = 1, 2, \dots, n. \quad (4.28)$$

<sup>4379</sup> or

$$y_{a,i} = \hat{b}x_{a,i} + e_i, \quad i = 1, 2, \dots, n, \quad (4.29)$$

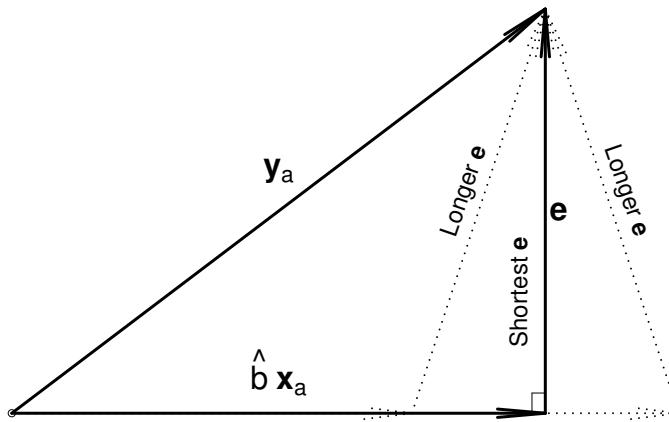
<sup>1</sup> SSE can be written as

$$\text{SSE} = |\mathbf{e}|^2 = |\mathbf{y}_a - \hat{b}\mathbf{x}_a|^2.$$

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  $\mathbf{e} \cdot \mathbf{x}_a = 0$ .

4380 or

$$\mathbf{y}_a = \hat{b} \mathbf{x}_a + \mathbf{e}. \quad (4.30)$$

4381 The dot product of both sides of Eq. (4.29) with  $\mathbf{x}_a$  yields

$$\hat{b} = \frac{\mathbf{y}_a \cdot \mathbf{x}_a - \mathbf{e} \cdot \mathbf{x}_a}{\mathbf{x}_a \cdot \mathbf{x}_a}. \quad (4.31)$$

4382 With the orthogonality condition Eq. (4.27)  $\mathbf{e} \cdot \mathbf{x}_a = 0$ , Eq. (4.31) is reduced to the  
4383 least square estimate of  $b$ : Eq. (4.21).

4384 In summary, the estimates  $\hat{a}$  and  $\hat{b}$  are derived by requiring (i) zero mean of the  
4385 residuals, and (ii) minimum sum of the residual squares.

4386 As shown in the computer code for Fig. 4.1, the R command for the simple linear  
4387 regression is `lm(y ~ x)`. The TLR dataset yields the estimate of the intercept  
4388  $\hat{a} = 33.476216$  and slope  $\hat{b} = -0.006982$ :

```
4389 lm(y ~ x)
4390 #(Intercept)      x
4391 #33.476216     -0.006982
```

4392 The estimated linear model is thus

$$y = 33.476216 - 0.006982x, \quad (4.32)$$

4393 OR

$$\text{Temp } [{}^{\circ}\text{C}] = 33.476216 - 0.006982 \times \text{Elevation } [\text{m}]. \quad (4.33)$$

4394 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]

```

4395

4396 The 24 residuals  $e_i, i = 1, 2, \dots, 24$  are

```

reg = lm(y ~ x)
round(reg$residuals, digits = 5)
#0.25792 1.53427 -0.37129 -0.43697 -0.10542 0.43358
# -0.60335 0.12833 -0.64953 -0.19817 0.10302 -0.63880
# -0.63366 -0.61692 -0.13283 0.63012 0.51454 1.27623
# -0.06333 0.27628 -1.52923 0.39122 1.52971 -1.09572

```

4403 The mean of the 24 residuals is zero, demonstrated by an R code below

```

mean(reg$residuals)
#[1] 1.62043e-17

```

4406 The least square condition  $\mathbf{e} \cdot \mathbf{x}_a = 0$  is demonstrated by the following R code

```

xa = x - mean(x)
sum(xa*reg$residuals)
#[1] -2.83773e-13

```

4410 The unbiased MSE is

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2. \quad (4.34)$$

4411 Here, subtracting 2 is due to the two constraints of estimating  $a$  and  $b$ . The unbiased  
4412 MSE  $s^2$  the Colorado TLR regression can be computed using the following R code

```

sum((reg$residuals)^2)/(length(y) -2)
#[1] 0.6096193

```

4415 Thus, one can use the data to estimate  $\hat{b}$  first using (4.21), then  $\hat{a}$  using (4.20),  
4416 and finally  $\hat{\sigma}^2$ , also denoted by  $s^2$ , using (4.34).

4417 The above computer code is summarized as follows.

```

#R code for the Colorado TLR regression analysis
lm(y ~ x)
#(Intercept)           x
#33.476216      -0.006982

reg = lm(y ~ x)
round(reg$residuals, digits = 5)
mean(reg$residuals)
#[1] 1.62043e-17

xa = x - mean(x)
sum(xa*reg$residuals)
#[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))
# -1.85037170708594e-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)
#.6096193452251238 #unbiased MSE

```

4431

4432 Once again, we emphasize that two assumptions are used in the derivation of the  
 4433 formulas to estimate  $\hat{a}$  and  $\hat{b}$ :

- 4434 (a) The unbiased model assumption: The mean residual is zero  $\bar{e} = 0$ , and  
 4435 (b) The optimization assumption: The residual vector  $e$  is perpendicular to the  
 4436  $x$ -anomaly vector  $\mathbf{x}_a = x - \bar{x}$ .

4437 Under these two assumptions, the estimators  $\hat{a}$  and  $\hat{b}$  are highly sensitive to the  
 4438  $Y$  outliers, particularly the outlier data corresponding to the smallest and largest  
 4439  $x$  values. One outlier can completely change the  $\hat{a}$  and  $\hat{b}$  values, in agreement with  
 4440 our intuition. This is an end point problem often encountered in data analysis. To  
 4441 suppress the sensitivity, many robust regression methods have been developed and  
 4442 R packages are available, such as the Robust Regression, by UCLA institute for  
 4443 Digital Research & Education:

4444 <https://stats.idre.ucla.edu/r/dae/robust-regression/>

4445 Figure 4.2 can be generated by the following R code.

```

#R plot Fig. 4.2: Geometric derivation of the least squares
4446 par(mar=c(0.0,0.5,0.0,0.5))
4447 plot(0,0, xlim=c(0,5.2), ylim=c(0,2.2),
4448       axes = FALSE, xlab="", ylab="")
4449 arrows(0,0,4,0, angle=5, code=2, lwd=3, length=0.5)
4450 arrows(4,0,4,2, angle=5, code=2, lwd=3, length=0.5)
4451 arrows(0,0,4,2, angle=5, code=2, lwd=3, length=0.5)
4452 arrows(5,0,4,2, angle=7, code=2, lwd=2, lty=3, length=0.5)
4453 arrows(0,0,5,0, angle=7, code=2, lwd=2, lty=3, length=0.5)
4454 arrows(3,0,4,2, angle=7, code=2, lwd=2, lty=3, length=0.5)
4455 arrows(0,0,3,0, angle=7, code=2, lwd=2, lty=3, length=0.5)
4456 segments(3.9,0, 3.9, 0.1)
4457 segments(3.9, 0.1, 4.0, 0.1)
4458 text(2,0.2, expression(hat(b)^bold(x)[a]), cex=2)
4459 text(2,1.2, expression(bold(y)[a]), cex=2)
4460 text(4.1,1, expression(bold(e)), cex=2)

```

```

4462 text(3.8,0.6, expression(paste("Shortest",bold(e))),  
4463     cex=1.5, srt=90)  
4464 text(3.4,1.1, expression(paste("Longer",bold(e))),  
4465     cex=1.5, srt=71)  
4466 text(4.6,1.1, expression(paste("Longer",bold(e))),  
4467     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()

```

4468

#### 4.1.2.3 Relationship between slope and correlation

4469 The slope  $\hat{b}$  and correlation  $r_{xy}$  are related in the following way

$$\hat{b}|\mathbf{x}_a| = r_{xy}|\mathbf{y}_a|. \quad (4.35)$$

4470 This has two extreme cases:  $r_{xy} = 1$  and  $r_{xy} = 0$ .

4471 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

$$\hat{b}|\mathbf{x}_a| = |\mathbf{y}_a|. \quad (4.36)$$

4474 Case (b). Perfect noise with zero correlation  $r_{xy} = 0$ , which implies

$$\hat{b} = 0, \quad (4.37)$$

4475 i.e., the best prediction is the mean.

4476 Formula (4.35) can be derived as follows. The slope estimate formula (4.21) can  
4477 be re-written to provide geometric and science interpretations.

4478 (i) Geometric projection interpretation:

$$\hat{b} = \frac{\mathbf{y}_a \cdot (\mathbf{x}_a / |\mathbf{x}_a|)}{|\mathbf{x}_a|}. \quad (4.38)$$

4479 Since  $\mathbf{x}_a / |\mathbf{x}_a|$  is a unit vector in the direction of  $\mathbf{x}_a$ , the slope is the projec-  
4480 tion of the  $y$  anomaly data vector on the  $x$  anomaly data vector and then  
4481 normalized by the  $x$  anomaly data vector.

4482 Or we write the  $\hat{b}$  formula in the following way:

$$\hat{b} = \frac{\mathbf{y}_a}{|\mathbf{x}_a|} \cdot \frac{\mathbf{x}_a}{|\mathbf{x}_a|}. \quad (4.39)$$

4483 The first fraction on the right hand side is the  $y$  anomaly data vector  
4484 normalized by  $|\mathbf{x}_a|$ . Thus, the slope is the projection of this normalized  $y$   
4485 anomaly data vector onto the unit  $x$  anomaly data vector.

4486 (ii) Correlation interpretation:

$$\hat{b} = r_{xy} \frac{|\mathbf{y}_a|}{|\mathbf{x}_a|} \quad (4.40)$$

4487 where  $r_{xy}$  is the correlation coefficient, or simply called the correlation,  
4488 defined by

$$r_{xy} = \frac{\mathbf{y}_a \cdot \mathbf{x}_a}{|\mathbf{x}_a| |\mathbf{y}_a|} \quad (4.41)$$

4489 The slope depends directly on the correlation coefficient between the  $x$  and  
4490  $y$  anomaly data. The correlation is scaled by the ratio of the length of the  
4491  $y$  anomaly data vector to the length of the  $x$  anomaly data vector.

4492 Below is the computer code corresponding to the two interpretations for the  
4493 Colorado July TLR example.

```
4494 #R code for estimating regression slope b
4495
4496 #Method 1: Using vector projection
4497 xa = x - mean(x) #Anomaly the x data vector
4498 nxa = sqrt(sum(xa^2)) #Norm of the anomaly data vector
4499 ya = y - mean(y)
4500 nya=sqrt(sum(ya^2))
4501 sum(ya*(xa/nxa))/nxa #Compute b
4502 #[1] -0.006981885 #This is an estimate for b
4503
4504 #Method 2: Using correlation
4505 corxy=cor(xa, ya) #Compute the correlation between xa and ya
4506 corxy
4507 #[1] -0.9813858 #Very high correlation
```

```
4508 corxy*nya/nxa #Compute b
4509 #[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
```

4510

#### 4.1.2.4 Percentage of variance explained: $R^2 \times 100\%$

4511 The unbiased variance of the model-predicted data  $\hat{y}_i(i = 1, 2, \dots, n)$  is defined as

$$\text{MV} = \frac{\sum_{i=1}^n [(\hat{y}_i - \bar{\hat{y}})]^2}{n - 1}, \quad (4.42)$$

4512 where  $\bar{\hat{y}}$  is the mean of  $\hat{y}_i(i = 1, 2, \dots, n)$ .

4513 The unbiased variance of the  $Y$  data

$$\text{YV} = \frac{\sum_{i=1}^n [(y_i - \bar{y})]^2}{n - 1}, \quad (4.43)$$

4514 where  $\bar{y}$  is the mean of the original station temperature data  $y_i(i = 1, 2, \dots, n)$ .

4515 We wish to measure how good the model is in terms of variance, and thus define  
4516 the ratio of MV to YV, which is named  $R^2$

$$R^2 = \frac{\text{MV}}{\text{YV}}, \quad (4.44)$$

4517 The value  $R^2 \times 100\%$  indicates the percentage of variance explained by the linear  
4518 model. A larger  $R^2$  value suggests a better model.

4520 The variance of the model-predicted data can be re-written as

$$\begin{aligned}
 \text{MV} &= \frac{\sum_{i=1}^n \left[ (\hat{a} + \hat{b}x_i) - (\hat{a} + \hat{b}\bar{x}) \right]^2}{n-1} \\
 &= \hat{b}^2 \frac{\sum_{i=1}^n [x_i - \bar{x}]^2}{n-1} \\
 &= \left[ r_{xy} \frac{|\mathbf{y}_a|}{|\mathbf{x}_a|} \right]^2 \frac{|\mathbf{x}_a|^2}{n-1} \\
 &= r_{xy}^2 \frac{|\mathbf{y}_a|^2}{n-1} \\
 &= r_{xy}^2 \times \text{YV}.
 \end{aligned} \tag{4.45}$$

4521 Therefore,

$$R^2 = r_{xy}^2. \tag{4.46}$$

4522 The  $R^2$  value is equal to the square of the correlation coefficient computed from  
4523 the data  $(x_i, y_i) (i = 1, 2, \dots, n)$  with  $0 \leq R^2 \leq 1$ .

4524 For the Colorado TLR data used earlier, we have  $r_{xy} = -0.9813858$ , and  $R^2 =$   
4525  $0.9631181$ .

4526 The variances MV, YV, and  $R^2$  can be computed by the following computer code  
4527 in multiple ways.

```

4528 #R code for computing MV
4529 var(reg$fitted.values)
4530 #[1] 15.22721
4531 #Or another way
4532 yhat = reg$fitted.values
4533 var(yhat)
4534 #[1] 15.22721
4535 #Or still another way
4536 n = 24
4537 sum((yhat - mean(yhat))^2)/(n-1)
4538 #[1] 15.22721
4539
4540 #R code for computing YV
4541 sum((y - mean(y))^2)/(n-1)
4542 # [1] 15.81033
4543 #Or another way
4544 var(y)
4545 #[1] 15.81033
4546
4547 #R code for computing R-squared value
4548 var(reg$fitted.values)/var(y)
4549 #[1] 0.9631181 #This is the R-squared value
4550
4551 cor(x,y)
4552 #[1] -0.9813858
4553 (cor(x,y))^2
4554 #[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

```

4555

#### 4.1.2.5 The regression estimates by the least square or the maximum likelihood

4558 A more complicated and commonly used derivation is based on the minimization  
 4559 of the sum of squared errors (SSE)

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n [y_i - (\hat{a} + \hat{b}x_i)]^2. \quad (4.47)$$

4560 This approach can be found in most statistics books and from the Internet. The  
 4561 terminology of “least square” regression comes from this minimization principle.

4562 This minimization is with respect to  $\hat{a}$  and  $\hat{b}$ . The minimization condition leads  
 4563 to two linear equations that determine  $\hat{a}$  and  $\hat{b}$ . The “least square” minimization  
 4564 condition for  $\hat{b}$  is equivalent to the orthogonality condition (4.27). Geometrically,  
 4565 the minimum distance of a point to a line is defined as the length of the line  
 4566 segment that is orthogonal to the line and connects the point to the line, that is, it

4567 is the minimum distance between the point and the line (See Fig. 4.2). Thus, the  
 4568 orthogonality condition and minimization condition are equivalent.

4569 Another commonly used method to estimate the parameters is to maximize a  
 4570 likelihood function defined as

$$L(a, b, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp(-(y_i - (a + bx_i))^2 / (2\sigma^2)). \quad (4.48)$$

4571 The solution for the maximum  $L(a, b, \sigma)$  yields the estimate of  $\hat{a}$ ,  $\hat{b}$ , and  $\hat{\sigma}^2$ . The  
 4572 result is exactly the same as the ones obtained by the method of least squares or  
 4573 perpendicular projection.

4574 However, the maximum likelihood approach shown here explicitly assumes the  
 4575 normal distribution of the error term  $\epsilon_i$  and the response variable  $Y_i$ . For the least  
 4576 squares approach, the assumption of normal distribution is needed only when for  
 4577 the inference of the parameters  $\hat{a}$ ,  $\hat{b}$ , and  $\hat{\sigma}^2$ , but is not needed to estimate them.

#### 4578 4.1.3 Statistics of slope and intercept: Distributions, confidence 4579 intervals, and inference

---

##### 4580 4.1.3.1 Mean and variance of the slope

4582 The normal distribution of  $a$  and  $b$  is now assumed in this section. This assumption  
 4583 was not needed earlier to estimate  $\hat{a}$  and  $\hat{b}$  and to compute R-squared.

4584 Different datasets will yield different estimates of  $\hat{a}$ ,  $\hat{b}$ , and  $\hat{\sigma}^2$ , which form a dis-  
 4585 tribution corresponding to the random datasets. Thus, we may regard  $\hat{a}$ ,  $\hat{b}$ , and  $\hat{\sigma}^2$   
 4586 as random variables. Their expected values are  $a$ ,  $b$  and  $\sigma^2$  if the random datasets  
 4587 satisfy the linear model assumptions (linearity, constant variance, independent er-  
 4588 rors, and normal distribution). However, in practical applications, we just have one  
 4589 dataset, estimate  $\hat{a}$ ,  $\hat{b}$ , and  $\hat{\sigma}^2$  once, and then interpret the results.

4590 Following Eq. (4.21), instead of using data  $y_i$  to estimate the slope, we use the  
 4591 corresponding random variables  $Y_i$  to define the slope  $B$  as a random variable:

$$\begin{aligned} B &= \frac{\mathbf{Y}_a \cdot \mathbf{x}_a}{\mathbf{x}_a \cdot \mathbf{x}_a} \\ &= \frac{\sum_{i=1}^n (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) Y_i, \end{aligned} \quad (4.49)$$

4592 because

$$\sum_{i=1}^n (x_i - \bar{x}) = 0. \quad (4.50)$$

4593 The above shows that  $B$  is a linear combination of  $n$  terms  $Y_i$ . If  $Y_i \sim N(a + bx, \sigma^2)$ ,  
 4594 then  $B$  is also normally distributed when  $n$  is large, because of the Central Limit  
 4595 Theorem.

4596 The expected value of  $B$  is  $b$ , as shown below:

$$\begin{aligned}
 E[B] &= \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) E[Y_i] \\
 &= \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) (a + bx_i) \\
 &= b \times \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) x_i \\
 &= b \times \left( \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \\
 &= b.
 \end{aligned} \tag{4.51}$$

4597 The variance of  $B$  is

$$\text{Var}(B) = \frac{\sigma^2}{\|\mathbf{x}\|_a^2}. \tag{4.52}$$

4598 This can be shown as follows:

$$\begin{aligned}
 \text{Var}[B] &= \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^2 \text{Var}(Y_i) \\
 &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
 &= \frac{\sigma^2}{\|\mathbf{x}\|_a^2}.
 \end{aligned} \tag{4.53}$$

4599 Statistics literature often denotes

$$S_{xx} = \|\mathbf{x}\|_a^2. \tag{4.54}$$

4600 Hence,

$$\text{Var}(B) = \frac{\sigma^2}{S_{xx}}. \tag{4.55}$$

4601 When  $n$  goes to infinity,  $S_{xx}$  also goes to infinity. Thus,

$$\lim_{n \rightarrow \infty} \text{Var}(B) = 0. \tag{4.56}$$

4602 This means that more data points yield a better estimate for the slope, supporting  
4603 our intuition.

#### 4604 4.1.3.2 Mean and variance of the intercept

4605 Instead of estimating the intercept using data  $y_i$  following formula (4.20), we define  
4606 the intercept as a random variable using the same formula but with random variable  
4607  $Y$ :

$$A = \bar{Y} - B\bar{x}. \tag{4.57}$$

<sup>4608</sup> Below we show that  $E[A] = a$ :

$$\begin{aligned} E[A] &= E[\bar{Y}] - E[B]\bar{x} \\ &= \frac{\sum_{i=1}^n Y_i}{n} - b \frac{\sum_{i=1}^n x_i}{n} \\ &= \frac{\sum_{i=1}^n (Y_i - bx_i)}{n} \\ &= \frac{\sum_{i=1}^n a}{n} \\ &= a. \end{aligned} \tag{4.58}$$

<sup>4609</sup> The variance of  $A$  is equal to

$$\text{Var}[A] = \frac{\sigma^2}{n} \left( 1 + \frac{\bar{x}^2}{\sigma_x^2} \right), \tag{4.59}$$

<sup>4610</sup> where

$$\sigma_x^2 = \frac{S_{xx}}{n} \tag{4.60}$$

<sup>4611</sup> may be considered the estimated variance of the  $x_i$  data. The result can be derived as follows:

$$\begin{aligned} \text{Var}[A] &= \text{Var}[\bar{Y}] + \text{Var}[B]\bar{x}^2 \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2}{S_{xx}} \bar{x}^2 \\ &= \frac{\sigma^2}{n} \left( 1 + \frac{\bar{x}^2}{S_{xx}/n} \right). \end{aligned} \tag{4.61}$$

<sup>4613</sup> This expression implies that as  $n$  goes to infinity, the variance  $\text{Var}[A]$  also goes to zero. Thus, the intercept estimate is better when having more data points.

#### <sup>4615</sup> 4.1.3.3 Confidence intervals of slope and intercept

<sup>4616</sup> The  $(1 - \alpha) \times 100\%$  confidence interval of the slope is

$$\left( \hat{b} - t_{\alpha/2, n-2} \frac{s}{\sqrt{S_{xx}}}, \hat{b} + t_{\alpha/2, n-2} \frac{s}{\sqrt{S_{xx}}} \right), \tag{4.62}$$

<sup>4617</sup> where

$$s = \sqrt{\frac{SSE}{n-2}} \tag{4.63}$$

<sup>4618</sup> is the unbiased estimator of  $\sigma$ , and the degrees of freedom (dof) of the t-distribution is  $n - 2$ .

<sup>4620</sup> If the confidence interval does not include zero, then the slope is significantly different from zero at the  $\alpha \times 100\%$  significance level. The Colorado TLR example has  $\hat{b} = -0.006982$ . For the 95% confidence interval, the R's quantile function `qt` of the t-distribution can produce the quantile interval corresponding to probabilities 0.025 and 0.975:

```
4625 qt(c(.025, .975), df=22)
4626 #[1] -2.073873 2.073873
```

4627 The estimated standard error for  $\hat{b}$  is

$$\text{SE}[\hat{b}] = \frac{s}{\sqrt{S_{xx}}} = 0.0002913 \quad (4.64)$$

4628 based on the R output from `summary(reg)`. Thus, the confidence interval for  $\hat{b}$  is

$$\begin{aligned} & (-0.0069818 - 2.073873 \times 0.0002913, -0.0069818 + 2.073873 \times 0.0002913) \\ & = (-0.0076, -0.0064). \end{aligned} \quad (4.65)$$

4629 This confidence interval does not include zero, and hence the slope is significantly  
4630 less zero at the 5% confidence level. The 95% confidence interval for TLR is (6.4,  
4631 7.6) °C/km based on the Colorado data used in this chapter.

4632 The  $(1 - \alpha) \times 100\%$  confidence interval of the intercept is

$$\left( \hat{a} - t_{\alpha/2, n-2} s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}, \hat{a} + t_{\alpha/2, n-2} s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} \right). \quad (4.66)$$

4633 The R output from `summary(reg)` gives the standard error of  $\hat{a}$

$$\text{SE}[\hat{a}] = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} = 0.5460279 \quad (4.67)$$

4634 Thus, the confidence interval for  $\hat{a}$  is

$$\begin{aligned} & (33.4762157 - 2.073873 \times 0.5460279, 33.4762157 + 2.073873 \times 0.5460279) \\ & = (32.3438, 34.6086). \end{aligned} \quad (4.68)$$

4635 The mathematical expressions of the confidence interval can be derived as follows.  
4636 Formulas (4.51), (4.55) (4.58) and (4.59) imply that

$$A \sim N \left( a, \frac{\sigma^2}{n} \left( 1 + \frac{\bar{x}^2}{\sigma_x^2} \right) \right), \quad (4.69)$$

$$B \sim N \left( b, \frac{\sigma^2}{S_{xx}} \right) \quad (4.70)$$

$$\sigma_x^2 = S_{xx}/n \quad (4.71)$$

4637 for a given  $\sigma^2$ .

4638 When  $\sigma^2$  is unknown and is to be estimated by  $s^2$ , then  $A$  and  $B$  follow a t-distribution.  
4639 To be exact, the t-statistics with  $n - 2$  dof can be defined as

$$T_a = \frac{\hat{a} - a}{s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}}, \quad (4.72)$$

$$T_b = \frac{\hat{b} - b}{s / \sqrt{S_{xx}}}, \quad (4.73)$$

4640 where  $a$  and  $b$  are given, and  $\hat{a}$  and  $\hat{b}$  are estimated from data. The confidence

4641 interval  $(-t_{\alpha/2,n-2}, t_{\alpha/2,n-2})$  for  $T_a$  and  $T_b$  leads to the confidence intervals for  $\hat{a}$   
4642 and  $\hat{b}$ . Thus, these two statistics  $T_a$  and  $T_b$  can be used for a hypothesis test.

#### 4643 4.1.3.4 Hypothesis test for the slope using the t-test

4644 We wish to test the hypothesis whether the slope is equal to a given constant  $\beta$ .  
4645 The null and alternative hypotheses are

$$H_0 : b = \beta, \quad (4.74)$$

$$H_1 : b \neq \beta. \quad (4.75)$$

4646 We define a t-statistic

$$T_b = \frac{\hat{b} - \beta}{s/\sqrt{S_{xx}}}. \quad (4.76)$$

4647 At the  $\alpha \times 100\%$  significance level, we reject  $H_0$  if

$$|T_b| > t_{1-\alpha/2,n-2}. \quad (4.77)$$

4648 For the Colorado July TLR example, if we want to test whether the TLR is equal  
4649 to  $7.3 \text{ }^{\circ}\text{C}/\text{km}$ , we calculate

$$T_b = \frac{-0.0069818 - (-0.0073)}{0.0002913} = 1.092345. \quad (4.78)$$

4650 The quantile

$$t_{0.975,22} = 2.073873 \quad (4.79)$$

4651 Thus,  $|T_b| < t_{0.975,22}$ , and hence the null hypothesis is not rejected at the 5%  
4652 significance level. Namely, the obtained TLR  $7.0 \text{ }^{\circ}\text{C}/\text{km}$  is not significantly different  
4653 from the given value  $7.3 \text{ }^{\circ}\text{C}/\text{km}$ .

#### 4654 4.1.3.5 Confidence intervals of prediction: The fitted model and 4655 the response variable

4656 For an assigned point  $x^*$ , which is not on the data points  $x_i$ , the fitted model  
4657 prediction is

$$\hat{Y} = \hat{a} + \hat{b}x^*. \quad (4.80)$$

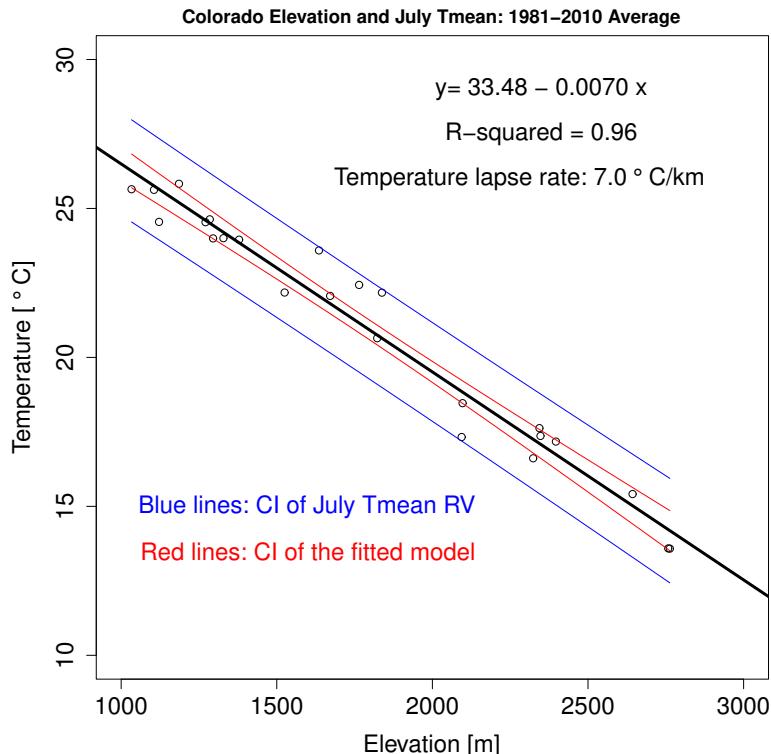
4658 The variance of this model prediction is

$$\text{Var}[\hat{a} + \hat{b}x^*] = \sigma^2 \left[ \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right]. \quad (4.81)$$

4659 Thus, the  $(1 - \alpha) \times 100\%$  confidence interval for  $\hat{a} + \hat{b}x^*$  is given by

$$\hat{a} + \hat{b}x^* \pm t_{1-\alpha/2,n-2}s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}. \quad (4.82)$$

This interval can be computed for any  $x = x^*$  and is shown by the red lines in Fig. 4.3. It means that 95% of the chance for the model prediction to be between the red lines. Note that the red lines diverge as  $x^*$  deviates from  $\bar{x}$ . This implies that 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).

The error of the prediction at  $x^*$  is

$$e^* = Y - \hat{Y} = Y - (\hat{a} + \hat{b}x^*) \quad (4.83)$$

is normally distributed if all the assumptions for the linear model are satisfied. The error's mean is zero and variance is

$$\begin{aligned} \text{Var}[e^*] &= \text{Var}[Y - \hat{a} + \hat{b}x^*] \\ &= \text{Var}[Y] + \text{Var}[\hat{a} + \hat{b}x^*] \\ &= \sigma^2 + \sigma^2 \left[ \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right]. \end{aligned} \quad (4.84)$$

This result implies that the confidence interval for the response variable  $Y$  at an

4668 assigned  $x^*$  is

$$\hat{a} + \hat{b}x^* \pm t_{1-\alpha/2, n-2}s\sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}. \quad (4.85)$$

4669 The confidence interval is shown by the blue lines in Fig. 2.2. It means that 95% of  
 4670 the chance for the  $Y$  values to be in between the blue lines. As expected, almost all  
 4671 the Colorado TLR observed data lie between the blue lines. Only 5% of chance a  
 4672 data point would lie outside the blue lines. According to Eq. (4.85), the blue lines  
 4673 also diverge, but very slowly because

$$\frac{(x^* - \bar{x})^2}{S_{xx}} < 0.1310$$

4674 is much smaller than  $1 + \frac{1}{n} = 1.0417$ .

4675 Figure 4.3 can be produced by the following computer code.

```

4676 #R plot Fig. 4.3: Confidence intervals of a regression model
4677 setwd("/Users/sshen/climstats")
4678 #Confidence interval of the linear model
4679 x1 = seq(max(x), min(x), len=100)
4680 n = 24
4681 xbar = mean(x)
4682 reg = lm(y ~ x)
4683 SSE = sum((reg$residuals)^2)
4684 s_squared = SSE/(length(y)-2)
4685 s = sqrt(s_squared)
4686 modTLR = 33.476216 + -0.006982*x1
4687 xbar = mean(x)
4688 Sxx = sum((x-xbar)^2)
4689 CIupperModel= modTLR +
4690   qt(.975, df=n-2)*s*sqrt((1/n)+(x1-xbar)^2/Sxx)
4691 CIlowerModel= modTLR -
4692   qt(.975, df=n-2)*s*sqrt((1/n)+(x1-xbar)^2/Sxx)
4693 CIupperResponse= modTLR +
4694   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x1-xbar)^2/Sxx)
4695 CIlowerResponse= modTLR -
4696   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x1-xbar)^2/Sxx)
4697
4698 setEPS() #Plot the figure and save the file
4699 postscript("fig0403.eps", height = 8, width = 8)
4700 par(mar=c(4.5,4.5,2.0,0.5))
4701 plot(x,y,
4702       ylim=c(10,30), xlim=c(1000,3000),
4703       xlab="Elevation [m]",
4704       ylab=bquote("Temperature[" ~ degree ~ "C]"),
4705       main="Colorado[Elevation and July Tmean: 1981-2010] Average",
4706       cex.lab=1.5, cex.axis=1.5)
4707 lines(x1,CIupperModel,type="l",col='red')
4708 lines(x1,CIlowerModel,type="l",col='red')
4709 lines(x1,CIupperResponse,type="l",col='blue')
4710 lines(x1,CIlowerResponse,type="l",col='blue')
4711 abline(reg,lwd=3)
4712 text(2280, 26,
4713 bquote("Temperature[lapse rate: 7.0" ~ degree ~ "C/km"]),
4714 cex=1.5)

```

```

4715 text(2350, 27.5, "R-squared=0.96", cex=1.5)
4716 text(2350, 29, "y=33.48-0.0070x", cex=1.5)
4717 text(1600, 15, "Blue lines: CI of July Tmean RV",
4718   col="blue", cex=1.5)
4719 text(1600, 13.5, "Red lines: CI of the fitted model",
4720   col="red", cex=1.5)
4721 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\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()

```

4722

#### 4.1.3.6 When the assumptions for the linear regression model are violated

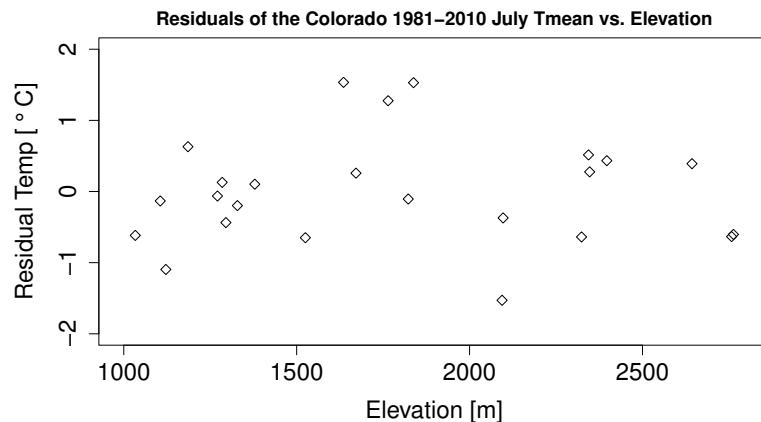
4723 The computing and plotting of a linear regression is quite straightforward, and  
 4724 have been used widely in data analysis. However, users often are unaware of the  
 4725 assumptions behind the computing and results. You may wish to know what are  
 4726 the assumptions and how they are violated.

4729 (a) **Review of the assumptions of linear regression**

4730 The linear regression model and its inference have four major assumptions:

- 4731 (i) Approximate linearity: The data support an approximate linear relationship  
4732 between  $x$  and  $Y$ ;
- 4733 (ii) Normality: Both the model error term  $\epsilon$  and  $Y$  are normally distributed;
- 4734 (iii) Constant variance:  $\epsilon$  and  $Y$  have variance  $\sigma^2$ , which does not change with  
4735 respect to  $x$ ;
- 4736 (iv) Independence of the error term:  $\text{Cor}(\epsilon_i, \epsilon_j) = \delta_{ij}$ .

4737 Although the derivation of the estimation formulas for  $\hat{a}$  and  $\hat{b}$  does not need any  
4738 of the above assumptions, the interpretation and inferences on  $\hat{a}$ ,  $\hat{b}$ ,  $\hat{a} + \hat{b}x$ , and  
4739 other quantities are often based on these assumptions. If one or more assumptions  
4740 are violated, the inference may not work. There are a variety of checking methods  
4741 and resolution methods in literature, such as using a nonlinear relationship (e.g., a  
4742 polynomial fitting model). R has a package `lmttest` to test linear regression models.  
4743 A residual scatter plot  $(x_i, e_i)$  may be a first good step to intuitively identify possible  
4744 violations of the assumptions (See Fig. 4.4). Ideally, the plot shows uniform noise for  
4745 a good model. Figure 4.4 shows the residuals of the Colorado temperature regression  
4746 against elevation, and seems showing this kind of noise without an obvious pattern.  
4747 We may safely conclude that the assumptions of linearity and constant variance are  
4748 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.

4749 To be rigorous on the statistical inferences about the results of a linear regres-  
4750 sion, a comprehensive residual analysis may be made to assess whether the studied  
4751 dataset satisfies the linear regression assumptions. If some assumptions are violated  
4752 for a given dataset, one may seek other methods to remediate the problems. One  
4753 way is to make a data transformation or to use a new regression function so that

4754 the assumptions are satisfied for the transformed data or the new functional model.  
 4755 Another way is to use a non-parametric method, which does not assume any distribution.  
 4756 The comprehensive residual analysis methods and the corresponding R  
 4757 and Python codes can be found online. This book is limited to the discussion of a  
 4758 few examples.

4759 Violation of independence assumption (iv) often occurs in cases of  $y$  being a time  
 4760 series. In this case there could be correlations from one time to later times. This  
 4761 serial correlation effect leads to fewer independent time intervals, i. e., the degrees of  
 4762 freedom (dof) can be reduced, and the statistical inference needs special attention.

```
4763 #R plot Fig. 4.4: Regression residuals
4764 reg = lm(y~x)
4765 setEPS() #Plot the figure and save the file
4766 postscript("fig0404.eps", height = 4.5, width = 8)
4767 par(mar=c(4.5,4.5,2.0,0.5))
4768 plot(x, reg$residuals, pch=5,
4769       ylim=c(-2,2), xlim=c(1000,2800),
4770       xlab="Elevation [m]",
4771       ylab=bquote("Residual[Temp]~degree~[C]"),
4772       main="Residuals of the Colorado July Tmean vs. Elevation",
4773       cex.lab=1.5, cex.axis=1.5, cex.main = 1.2)
4774 dev.off()
```

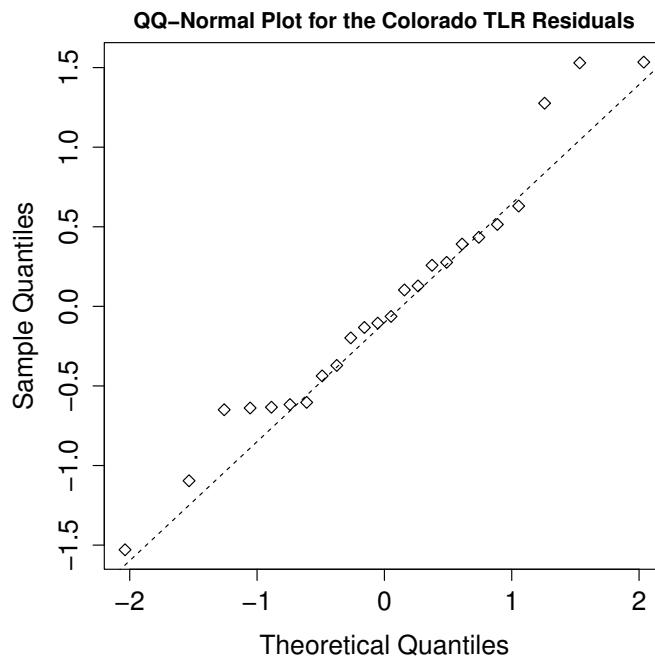
```
# Python plot Fig. 4.4: Regression residuals
reg1 = np.array(np.polyfit(x, y, 1))#regression
abline = reg1[1] + x*reg1[0]
# calculate residuals
r = y - abline
fig, ax = plt.subplots(figsize=(12,8))
ax.plot(x, r, 'kd')
ax.set_title("Residuals of the Colorado 1981-2010 July\nTmean vs. Elevation",fontweight = '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()
```

4775

### 4776 (b) Test for normality

4777 The normality and independence assumptions may not show in the residual scatter  
 4778 plot. You may use a Q-Q plot to test the normality, and Durbin-Watson test  
 4779 to check the independence. In the Colorado TLR example, the Q-Q plot is shown  
 4780 in Fig. 4.5. The residual data quantile points overlap fairly well with the theoretical  
 4781 line for the normal distribution. We may conclude that the data satisfy the  
 4782 normality assumption.

4783 However, the QQ-plot is only a subjective visual test. A quantitative test may  
 4784 be used, such as the Kolmogorov-Smirnov (KS) test or Shapiro-Wilk (SW) test.



**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.

4785 The KS-test statistic and its p-value for the Colorado TLR data are  $D = 0.20833$   
 4786 and  $p\text{-value} = 0.686$ . The null hypothesis is not rejected, i.e., the temperature  
 4787 residuals from the linear regression are not significantly different from the normal  
 4788 distribution. The normality assumption is valid.

4789 Figure 4.5 can generated by the following computer code.

```

4790 #R plot Fig. 4.5: Q-Q plot of quantiles
4791 reg = lm(y ~ x)
4792 setEPS() #Plot the figure and save the file
4793 postscript("fig0405.eps", height = 6, width = 6)
4794 par(mar=c(4.5,4.5,2.0,0.5))
4795 qqnorm(reg$residuals, pch=5,
4796   main="QQ-Normal Plot for the Colorado TLR Residuals",
4797   cex.lab = 1.4, cex.axis = 1.4)
4798 qqline(reg$residuals, lty=2)
4799 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()

```

4800

### 4801 (c) Serial correlation

4802 The R package `lmtest` has a function to test for the independence of the error  
 4803 term, using the Durbin-Watson (DW) test statistic defined as

$$DW = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}. \quad (4.86)$$

4804 The DW test is meant for a time series. We thus sort the Tmean data according to  
 4805 the ascending order of elevation, and then compute DW.

4806 When  $\text{cor}(e_i, e_j) = \sigma^2 \delta_{ij}$ , we have  $DW \approx 2$ . For the Colorado July TLR data,  
 4807  $DW = 2.3072$  and  $p-value = 0.7062$ , which can be computed by the following  
 4808 code:

```

4809 #R code for the DW-test for independence
4810 #install.packages("lmtest")
4811 library(lmtest)
4812 ElevTemp=cbind(x,y, 1:24)
4813 #Sort the data for ascending elevation
4814 ElevTemp=ElevTemp[order(x),]
4815 reg1=lm(ElevTemp[,2] ~ ElevTemp[,1])
4816 dwtest(reg1)
4817 #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

```

4818

4819 The large p-value 0.7062 implies that the null hypothesis of no serial correlation  
 4820 (i.e., independence) is not rejected, because DW 2.3072 is not too far away from 2,  
 4821 which indicates independence of the error terms. Some literatures use  $1.5 < DW <$   
 4822 2.5 to conclude independence without checking the p-value.

4823 When the independence assumption is violated, the data have serial correlation,  
 4824 which reduces dof and hence enlarges the confidence interval for the regression  
 4825 results and makes the results not as reliable.

#### 4826 (d) Non-parametric trend inference and heteroskedasty

4827 Even when some of the four assumptions for the linear regression model are  
 4828 invalid, one can still estimate  $\hat{a}$  and  $\hat{b}$  and the linear model  $\hat{a} + \hat{b}x$  in the same  
 4829 way and can still interpret the study subject using the domain knowledge. In this  
 4830 case, the statistical inference based on the normality assumption may not make  
 4831 any sense, and the confidence interval calculations based on the t-distribution are  
 4832 generally invalid. Now, you may use a non-parametric test to make an inference  
 4833 on the residuals. For example, a non-parametric test for the trend can be applied,  
 4834 such as Mann-Kendall (MK) trend test, and Theil-Sen's trend estimation and test.  
 4835 The trend, if presenting, can be linear or nonlinear. These tests are meant for time  
 4836 series. Thus, before applying these tests, the data should be sorted according to the  
 4837 ascending order of the explainable variable. For the Colorado TLR example, the  
 4838 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

```

4852

4853 The large p-value 0.6374 of the MK-test for residuals implies that the residuals  
 4854 have no trend, linear or nonlinear. The small p-value 2.261e-09 of the MK-test for  
 4855 the sorted temperature implies that the sorted July Tmean according to elevation  
 4856 has a significant trend.

4857 The linear trend is frequently used to explain the global average temperature. The  
 4858 temperature data with respect to time is clearly nonlinear. When the scatter plot  
 4859 shows a clear pattern, then one or more of the four assumptions are usually violated.  
 4860 When that happens, you may use different ways to remediate this situation. For  
 4861 example, when linearity is violated, the problem may be intrinsically nonlinear and  
 4862 a polynomial model may be a better fit than a linear model. When the constant  
 4863 variance is violated, known as heteroskedastic data in statistics literature, then the  
 4864 data standardization (i.e., the anomaly data being divided by the data standard  
 4865 deviation) or logarithmic transform for the positive-valued data, may be used to  
 4866 transform the data and to make the data homoskedastic, i.e., constant variance.

4867

## 4.2 Multiple linear regression

4868

4869 This section uses the method of matrices, which will be described in Chapter 5. If  
 4870 you are not familiar with matrices, you may come back to this section after reading  
 4871 the first two sections of Chapter 5.

4872

### 4.2.1 Calculating the Colorado TLR when taking location coordinates into account

4873

4874

4875 In the previous simple linear regression, the July temperature was assumed to  
 4876 linearly depend only on the vertical coordinate of the station: elevation. However,

4877 the temperature may also depend on the horizontal coordinates of the station:  
 4878 latitude and longitude. This subsection deals with this kind of problem of more  
 4879 than one explanatory variable. A multivariate linear regression model (or called  
 4880 multiple linear regression model) can be expressed as follows:

$$Y = b_0 + b_1x_1 + b_2x_2 + \cdots + b_mx_m + \epsilon, \quad (4.87)$$

4881 where  $x_1, x_2, \dots, x_m$  are  $m$  explanatory variables, which are non-random and de-  
 4882 terministic variables;  $\epsilon$  is the model error with zero mean and variance  $\sigma^2$  being a  
 4883 constant;  $Y$  is the response variable, which is random and has its variance equal to  
 4884  $\sigma^2$ , and expected value equal to

$$E[Y] = b_0 + b_1x_1 + b_2x_2 + \cdots + b_mx_m; \quad (4.88)$$

4885 and  $b_0, b_1, b_2, \dots, b_m$  are parameters to be estimated from data  
 4886  $(x_{ij}, y_j), i = 1, 2, \dots, m$ , and  $j = 1, 2, \dots, n$ .

4887 In the Colorado July TLR example, we have  $x_1$  as latitude,  $x_2$  as longitude,  $x_3$   
 4888 as elevation,  $Y$  as the July air temperature,  $m = 3$ , and  $n = 24$ . The latitude and  
 4889 longitude data are as follows:

```
4890 lat=c(
4891 39.9919, 38.4600, 39.2203, 38.8236, 39.2425, 37.6742,
4892 39.6261, 38.4775, 40.6147, 40.2600, 39.1653, 38.5258,
4893 37.7717, 38.0494, 38.0936, 38.0636, 37.1742, 38.4858,
4894 8.0392, 38.0858, 40.4883, 37.9492, 37.1786, 40.0583
4895 )
4896 lon=c(
4897 -105.2667, -105.2256, -105.2783, -102.3486, -107.9631, -106.3247,
4898 -106.0353, -102.7808, -105.1314, -103.8156, -108.7331, -106.9675,
4899 -107.1097, -102.1236, -102.6306, -103.2153, -105.9392, -107.8792,
4900 -103.6933, -106.1444, -106.8233, -107.8733, -104.4869, -102.2189
4901 )
```

4902 The elevation and temperature data are the same as those at the beginning of this  
 4903 chapter.

4904 The computer code for the multiple linear regression of three variables is as  
 4905 follows.

```
4906 #R code for the TLR multivariate linear regression
4907 elev = x; temp = y #The x and y data were entered earlier
4908 dat = cbind(lat, lon, elev, temp)
4909 datdf = data.frame(dat)
4910 datdf[1:2,] #Show the data of the first two stations
4911 #      lat      lon    elev    temp
4912 # 39.9919 -105.2667 1671.5 22.064
4913 # 38.4600 -105.2256 1635.6 23.591
4914
4915 #Multivariate linear regression
4916 reg=lm(temp ~ lat + lon + elev, data = datdf)
```

```

4917 summary(reg) #Display the regression results
4918 #           Estimate Std. Error t value Pr(>|t|)
4919 #(Intercept) 36.4399561 9.4355746 3.862 0.000971 ***
4920 #   lat        -0.4925051 0.1320096 -3.731 0.001319 **
4921 #   lon        -0.1630799 0.0889159 -1.834 0.081564 .
4922 #   elev       -0.0075693 0.0003298 -22.953 7.67e-16 ***
4923 #Residual standard error: 0.6176 on 20 degrees of freedom
4924 #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_)

```

4925

4926 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  
4927 estimated earlier using only one variable: elevation.  
4928

4929 The quantile of  $t_{0.975,20} = 2.085963$ . The dof is 20 now because four parameters  
4930 have been estimated. The confidence interval of TLR can be computed below

$$\begin{aligned}
& (-0.0075694 - 2.085963 \times 0.0003298, -0.0075694 + 2.085963 \times 0.0003298) \\
& = (-0.008257351, -0.006881449). \tag{4.89}
\end{aligned}$$

4931 The confidence interval for TLR at 95% confidence level is  $(6.9, 8.3)^{\circ}\text{C}/\text{km}$ .

4932 The  $R^2$  values is very large: 0.979, which means that the linear model data can  
4933 explain 98% of the variance of observed temperature data.

---

4934    **4.2.2 Formulas for estimating parameters in the multiple linear  
4935    regression**

---

4937    The  $n$  groups of data for the explainable  $x_1, x_2, \dots, x_m$  and response variable  $Y$   
4938    are as follows:

$$x_{ij}, i = 1, 2, \dots, m, \quad \text{and} \quad j = 1, 2, \dots, n; \quad (4.90)$$

4939    and

$$y_1, y_2, \dots, y_n \quad (4.91)$$

4940    The data and their corresponding regression coefficients are written in the matrix  
4941    form as follows:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (4.92)$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{m1} \\ 1 & x_{12} & x_{22} & \cdots & x_{m2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{mn} \end{bmatrix} \quad (4.93)$$

$$\hat{\mathbf{b}} = \begin{bmatrix} \hat{b}_0 \\ \hat{b}_1 \\ \hat{b}_2 \\ \vdots \\ \hat{b}_n \end{bmatrix} \quad (4.94)$$

4942    The data, linear model prediction, and the corresponding residuals (i.e., the pre-  
4943    diction errors of the linear model) can be written as follows:

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times (m+1)} \hat{\mathbf{b}}_{(m+1) \times 1} + \mathbf{e}_{n \times 1} \quad (4.95)$$

4944    where the residual vector is

$$\mathbf{e}_{n \times 1} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad (4.96)$$

4945    Multiplying Eq. (4.95) by  $\mathbf{X}^t$  from the left and enforcing

$$(\mathbf{X}^t)_{(m+1) \times n} \mathbf{e}_{n \times 1} = \mathbf{0}, \quad (4.97)$$

4946    you can obtain the estimate of the regression coefficients

$$\hat{\mathbf{b}}_{(m+1) \times 1} = [(\mathbf{X}^t)_{(m+1) \times n} \mathbf{X}_{n \times (m+1)}]^{-1} (\mathbf{X}^t)_{(m+1) \times n} \mathbf{y}_{n \times 1}. \quad (4.98)$$

Condition (4.97) means the each column vector of the  $\mathbf{X}_{n \times (m+1)}$  matrix is perpendicular to the residual vector  $\mathbf{e}_{n \times 1}$ . For the first column, the condition corresponds to the assumption of zero mean of the model error:

$$\sum_{i=1}^n e_i = 0. \quad (4.99)$$

For the remaining columns, the condition means that the residual vector is perpendicular to the data vector for each explanatory variable. This implies that the residual vectors Euclidean distance to the data vectors are minimized, i.e., the minimizing sum of squared errors (SSE). Therefore, Equation (2.83) is the least square estimate of the regression coefficients.

As illustrated in the previous subsection, to implement the R estimate of the regression coefficients, we put data  $\mathbf{X}$  and  $\mathbf{y}$  in a single matrix in the form of `datdf=data.frame(cbind(X,y))` with proper column names, such as  
`colnames(datdf) <- c('x1', 'x2', 'x3', 'y')`  
for the case of three explanatory variables. Then, use the R command to make the linear model estimate

```
reg=lm(y ~ x1 + x2 + x3, data = datdf)
```

Finally, `summary(reg)` outputs all the important regression results.

The R command `reg$` allows to display the specific result, such as regression coefficients

```
round(reg$coefficients, digits=5)
# (Intercept)      lat          lon          elev
#     36.43996    -0.49251    -0.16308    -0.00757
```

Thus, the multiple linear model for the Colorado July temperature is

$$\text{Temp}[^{\circ}\text{C}] = 36.43527 - 0.49255 \times \text{Lat} - 0.16314 \times \text{Lon} - 0.00757 \times \text{Elev [m]}. \quad (4.100)$$

The confidence interval for the linear model

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + \cdots + b_m x_m \quad (4.101)$$

at a given point  $\mathbf{x}^* = (1, x_1^*, x_2^*, \dots, x_m^*)$  is given by the following formula

$$\hat{\mathbf{b}}^t \mathbf{x}^* \pm t_{1-\alpha/2, n-m} s \sqrt{(\mathbf{x}^*)^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}^*}. \quad (4.102)$$

The confidence interval for the response variable  $Y$  at a given point

$$\mathbf{x}^* = (1, x_1^*, x_2^*, \dots, x_m^*)$$

is

$$\hat{\mathbf{b}}^t \mathbf{x}^* \pm t_{1-\alpha/2, n-m} s \sqrt{1 + (\mathbf{x}^*)^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}^*}, \quad (4.103)$$

where  $s$  is the same as that defined earlier for the simple linear regression:

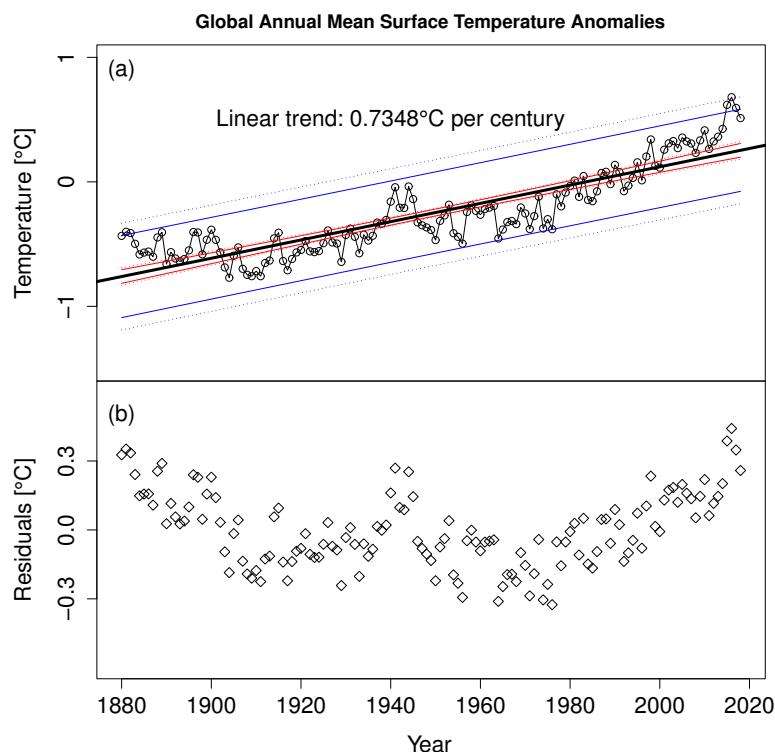
$$s = \sqrt{\frac{SSE}{n-2}} \quad (4.104)$$

If  $m = 1$ , the above two formulas are the same as those for the confidence intervals given by Eqs. (4.82) and (4.85) for the simple linear regression.

## 4.3 Nonlinear fittings using the multiple linear regression

### 4.3.1 Diagnostics of linear regression: An example of global temperature

When the data have strong nonlinearity, the scatter plot of residuals will show obvious patterns, such as that shown in Fig. 4.6 for the linear regression of the global average annual mean surface air temperature anomalies from 1880 to 2018 with respect to the 1971-2000 climatology. The linearity assumption of the simple 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.

The following computer code can plot Fig. 4.6 and make the diagnostics of the linear regression.

```
#R plot Fig. 4.6: Regression diagnostics
```

```

4990 setwd("/Users/sshenn/climstats")
4991 dtmean<-read.table(
4992   "data/aravg.ann.land_ocean.90S.90N.v5.0.0.201909.txt",
4993   header=F)
4994 dim(dtmean)
4995 #[1] 140 6
4996 x = dtmean[1:139,1]
4997 y = dtmean[1:139,2]
4998 reg = lm(y ~ x) #linear regression
4999 reg
5000 #(Intercept)      yrtime
5001 #-14.574841      0.007348
5002
5003 #Confidence interval of the linear model
5004 xbar = mean(x)
5005 SSE = sum((reg$residuals)^2)
5006 s_squared = SSE/(length(y)-2)
5007 s = sqrt(s_squared)
5008 modT = -14.574841 + 0.007348 *x
5009 xbar = mean(x)
5010 Sxx = sum((x-xbar)^2)
5011 n = length(y)
5012 CIupperModel= modT +
5013   qt(.975, df=n-2)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5014 CIlowerModel= modT -
5015   qt(.975, df=n-2)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5016 CIupperResponse= modT +
5017   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5018 CIlowerResponse= modT -
5019   qt(.975, df=n-2)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5020
5021 CIupperModelr= modT +
5022   qt(.975, df=5)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5023 CIlowerModelr= modT -
5024   qt(.975, df=5)*s*sqrt((1/n)+(x-xbar)^2/Sxx)
5025 CIupperResponser= modT +
5026   qt(.975, df=5)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5027 CIlowerResponser= modT -
5028   qt(.975, df=5)*s*sqrt(1+(1/n)+(x-xbar)^2/Sxx)
5029
5030 setEPS() #Plot the figure and save the file
5031 postscript("fig0406.eps", height = 8, width = 8)
5032 par(mfrow=c(2,1))
5033 par(mar=c(0,4.5,2.5,0.7))
5034 plot(x, y, ylim = c(-1.5, 1),
5035   type="o", xaxt="n", yaxt="n",
5036   cex.lab=1.4, cex.axis=1.4,
5037   xlab="Year", ylab=bquote("Temperature["~degree~"C]"),
5038   main="Global Annual Mean Surface Temperature Anomalies",
5039   cex.lab=1.4, cex.axis=1.4
5040 )
5041 axis(side = 2, at = c(-1.0, 0, 1.0), cex.axis = 1.4)
5042 abline(reg, col="black", lwd=3)
5043 lines(x,CIupperModel,type="l",col='red')
5044 lines(x,CIlowerModel,type="l",col='red')
5045 lines(x,CIupperResponse,type="l",col='blue')

```

```
5046 lines(x,CIlowerResponse ,type="l",col='blue')
5047
5048 lines(x,CIupperModelr ,type="l", lty = 3, col='red')
5049 lines(x,CIlowerModelr ,type="l", lty = 3, col='red')
5050 lines(x,CIupperResponser ,type="l",lty = 3, col='blue')
5051 lines(x,CIlowerResponser ,type="l",lty = 3, col='blue')
5052
5053 text(1940, 0.5,
5054   bquote("Linear trend:~0.7348`^degree`^"C_per_century"),
5055   col="black",cex=1.4)
5056 text(1880, 0.9, "(a)", cex=1.4)
5057 par(mar=c(4.5,4.5,0,0.7))
5058 plot(x, reg$residuals , ylim = c(-0.6,0.6),
5059   pch=5, cex.lab=1.4, cex.axis=1.4,
5060   yaxt = 'n', xlab="Year",
5061   ylab=bquote("Residuals`[^degree`^"C]"))
5062 axis(side = 2, at = c(-0.3, 0, 0.3), cex.axis = 1.4)
5063 text(1880, 0.5, "(b)", cex=1.4)
5064 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 Annual Mean Land and Ocean Surface\nTemperature Anomalies", fontweight = 'bold',
    size = 25, pad = 20)
ax[0].set_ylabel("Temperature [$\degree C]", size = 25, labelpad = 20)
ax[0].tick_params(length=6, width=2, labelsize=20);
ax[0].set_yticks(np.round(np.linspace(-1, 1, 5), 2))
ax[0].text(1880, 0.4,
           "Linear temp trend 0.7348 deg C per century",
           color = 'k', size = 20)
ax[0].text(1877, 0.6, "(a)", size = 20)
ax[0].axes.get_xaxis().set_visible(False)
ax[1].tick_params(length=6, width=2, labelsize=20);
ax[1].set_yticks(np.round(np.linspace(-0.5, 0.5, 5), 2))
ax[1].set_ylabel("Residuals [$\degree 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()

```

5065

Based on the visual check of Fig. 4.6(b), the constant variance assumption seems satisfied. The KS-test shows that the normality assumption is also satisfied.

```
5068 #Kolmogorov-Smirnov (KS) test for normality
5069 library(fitdistrplus)
5070 resi_mean = mean(reg$residuals)
5071 resi_sd = sd(reg$residuals)
5072 test_norm = rnorm(length(reg$residuals),
5073                      mean = 0, sd = 1)
```

```

5074 testvar = (reg$residuals - resi_mean)/resi_sd
5075 ks.test(testvar, test_norm)
5076 #D = 0.057554, p-value = 0.9754
5077 #The normality assumption is accepted
5078
5079 #Diagnostics on independence and normality
5080 # Durbin-Watson (DW) test for independence
5081 dwtest(reg)
5082 #DW = 0.45235, p-value < 2.2e-16
5083 #The independence assumption is rejected
5084
5085 #degrees of freedom and critical t values
5086 rho1 = acf(y)[[1]][2] #Auto-correlation function
5087 rho1 #[1] 0.9270817
5088 edof = (length(y) - 2)*(1 - rho1)/(1 + rho1)
5089 edof #[1] 5.183904 effective degrees of freedom
5090 qt(.975, df=137) #[1] 1.977431 critical t value
5091 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

```

5092

The DW-test shows that the independence assumption is violated. This implies the existence of serial correlation, which in turn implies a smaller dof, and hence larger  $t_{1-\alpha,\text{dof}}$ . Thus, the confidence intervals for model  $\hat{Y}$  are wider than the red lines in Fig. 4.6(a), and those for the  $Y$  values are wider than the blue lines in the same figure. Therefore, the regression results are less reliable. In this case,

you need to compute the effective dof (edof), which is less than  $n - 2$ . For the NOAAGlobalTemp time series, the one-year time lag serial correlation is  $\rho_1 = 0.9271$ , the edof under the assumption of an autoregression one process AR(1) is approximately equal to

$$\text{edof} = \frac{1 - \rho_1}{1 + \rho_1} \times \text{dof}. \quad (4.105)$$

Consequently, the effective dof is reduced to 5, compared to the non-serial correlation case dof 137. R can compute `qt(.975, df=137)` and yield 1.977431, and compute `qt(.975, df=5)` and yield 2.570582. Thus, the actual confidence interval at the 95% confidence level for the linear regression is about 25% wider. Namely, the dotted color lines are wider than the solid color lines.

The W pattern of the residuals shown in Fig. 4.6(b) implies that the linearity assumption is violated. To remediate the nonlinearity, we fit the data to a nonlinear model and hope the residuals do not show clear patterns. We will use the third-order polynomial to illustrate the procedure.

### 4.3.2 Fit a third order polynomial

Figure 4.7(a) shows the third order polynomial fitting to the NOAAGlobalTemp dataset:

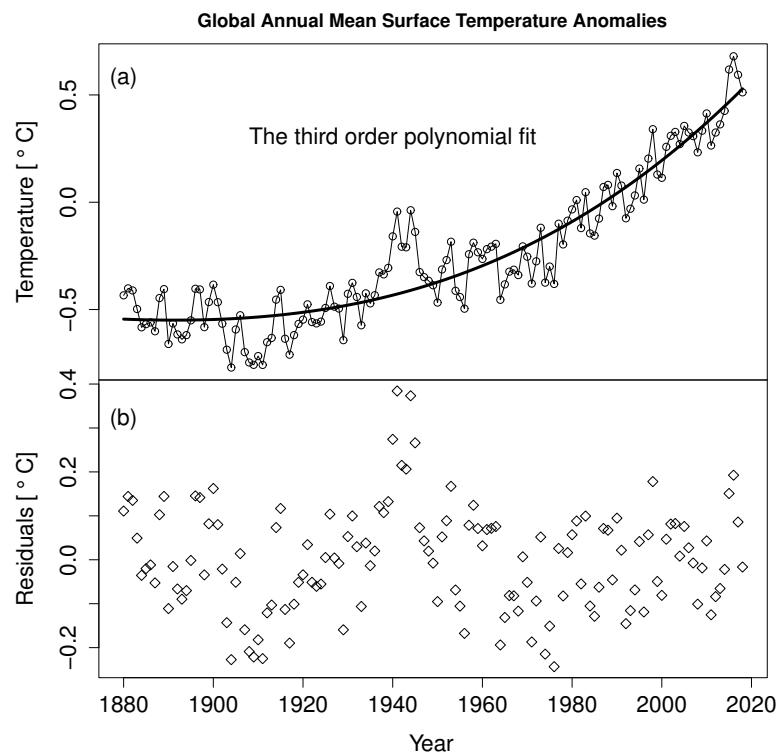
$$y = b_0 + b_1x + b_2x^2 + b_3x^3 + \epsilon. \quad (4.106)$$

This fitting can be estimated by the method of multiple linear regression with three variables

$$x_1 = x, x_2 = x^2, x_3 = x^3. \quad (4.107)$$

Then, a computer code for the multiple linear regression can be used to estimate the coefficients of the polynomial and plot Fig. 4.7.

```
5119 #R plot Fig. 4.7: Polynomial fitting
5120 x1=x
5121 x2=x1^2
5122 x3=x1^3
5123 dat3=data.frame(cbind(x1,x2,x3,y))
5124 reg3 = lm(y ~ x1 + x2 + x3, data=dat3)
5125 # simply use
5126 # reg3 = lm(y ~ x + I(x^2) + I(x^3))
5127 setEPS() #Plot the figure and save the file
5128 postscript("fig0407.eps", height = 8, width = 8)
5129 par(mfrow=c(2,1))
5130 par(mar=c(0,4.5,2.5,0.7))
5131 plot(x, y, type="o", xaxt="n",
5132       cex.lab=1.4, cex.axis=1.4, xlab="Year",
5133       ylab=bquote("Temperature["~degree~"C]"),
5134       main="Global~Annual~Mean~Surface~Temperature~Anomalies",
5135       cex.lab=1.4, cex.axis=1.4)
5136 lines(x, predict(reg3), col="black", lwd=3)
5137 reg3
5138 #(Intercept)           x1                  x2                  x3
```



**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.

```

5139 #-1.426e+03    2.333e+00   -1.271e-03    2.308e-07
5140 text(1940, 0.3,
5141     "The third order polynomial fit",
5142     col="black", cex=1.4)
5143 text(1880, 0.58, "(a)", cex=1.4)
5144 par(mar=c(4.5,4.5,0,0.7))
5145 plot(x1, reg3$residuals,
5146     pch=5, cex.lab=1.4, cex.axis=1.4,
5147     xlab="Year", ylab=bquote("Residuals["~degree~"C]"))
5148 text(1880, 0.32, "(b)", cex=1.4)
5149 dev.off()

```

```

os.chdir("/Users/sshen/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()

```

5150

5151 R has another simpler command for the third order polynomial regression:

5152 `lm(y ~ poly(x, 3, raw=TRUE))`

5153 where, `raw=TRUE` means using the raw polynomial form written like formula (4.106).  
 5154 Another option is `raw=FALSE`, which means that the data are fitted to an orthogonal  
 5155 polynomial.

5156 Comparison of Figs. 4.6(b) and 4.7(b) shows that the W-shape pattern of the  
 5157 residuals in the third-order polynomial fitting is weaker than that for the linear  
 5158 regression. Nonetheless, the W-shape pattern is still clear. Figure 4.7(b) visually  
 5159 suggests that the constant variance assumption is satisfied. The KS-test shows  
 5160 that the normality assumption is also satisfied. However, the DW-test shows the  
 5161 existence of serial correlation. The computing of these tests are left as exercise  
 5162 problems.

5163 You can try to fit a ninth order orthogonal polynomial. This can eliminate the W-

5164 shape nonlinear pattern of the residuals. You can also try to fit other polynomials  
5165 or functions.

## 5166 4.4 Linear Regression by Weighted Least Squares

5167

5168 The usual regression minimizes the mean square error:

$$\epsilon^2 = \sum_i (y_i - \hat{y}_i)^2. \quad (4.108)$$

5169 However, sometimes the square error may have different importance for different  
5170 observations  $i$ . The importance is quantified by a weight  $w_i$ . Thus, the weight least  
5171 square error is defined as

$$\varepsilon^2 = \sum_i w_i (y_i - \hat{y}_i)^2. \quad (4.109)$$

5172 The weights  $w_i$  are often required to be non-negative and sometimes to be normalized  
5173

$$\sum_i w_i = 1, \quad w_i \geq 0, \quad i = 1, 2, 3, \dots \quad (4.110)$$

5174 As an example of weights, the gridded global climate data are often defined on  
5175 a uniform latitude-longitude grid. A grid box covers more area over the equatorial  
5176 region than a polar region. From the spherical coordinates, you can find that the  
5177 area-weight is  $\cos \phi_i$ , where  $\phi_i$  is the latitude of the grid box's centroid.

5178 The weighted least squares regression can be converted into the usual least square  
5179 regression for the weighted field and data:

$$U_i = \sqrt{w_i} y_i, \quad \hat{U}_i = \sqrt{w_i} \hat{y}_i, \quad i = 1, 2, 3, \dots \quad (4.111)$$

5180 Here,  $\sqrt{w_i}$  are often referred to as weight-factor. In climate science,  $\sqrt{w_i} = \sqrt{\cos \phi_i}$   
5181 are called area-factor. In terms of the weighted data, the weighted least squares  
5182 becomes the usual least squares:

$$\varepsilon^2 = \sum_i (U_i - \hat{U}_i)^2. \quad (4.112)$$

5183 When considering 3D data in space, we may need to introduce the volume-factor  
5184  $\sqrt{V_i}$ . When heat capacity, mass-density and other physical or chemical properties  
5185 in addition to geometry need to be considered, then a more complex weight should  
5186 be used. For more examples, see Bui et al. (2023), Lafarga et al. (2023)), and Shen  
5187 and North (2023).

5188

## 4.5 Chapter summary

5189

5190 This chapter has introduced three regression methods commonly used in climate  
5191 science: (a) simple linear regression of a single variate, (b) multiple linear regression,  
5192 and (c) nonlinear fitting. For (a) and (b) we used the data of surface air temperature  
5193 and elevation of the 24 USHCN stations in the state of Colorado, USA. For (c) we  
5194 used the NOAAGlobalTemp's global average annual mean temperature data from  
5195 1880 to 2018.

5196 A linear regression model has four fundamental assumptions:

- 5197 (i) Approximate linearity between  $x$  and  $Y$ ,  
5198 (ii) Normal distribution of  $\epsilon$  and  $Y$ ,  
5199 (iii) Constant variance of  $\epsilon$  and  $Y$ , and  
5200 (iv) Independence of error terms:  $\text{Cor}(\epsilon_i, \epsilon_j) = \sigma^2 \delta_{ij}$ .

5201 Assumptions (i) and (iii) can be verified by visually examining the scatter plot  
5202 of residuals. Assumption (ii) may be verified by the KS-test. Assumption (iv) can  
5203 be verified by the DW-test. When one or more assumptions are violated, you may  
5204 improve your model, such as using a nonlinear model, or transforming the data.

## References and Further Readings

- 5206 [1] Cordova, M., R. Celleri, C.J. Shellito, J. Orellana-Alvear, A. Abril, and G.  
 5207 Carrillo-Rojas, 2016: Near-surface air temperature lapse rate over complex ter-  
 5208 rain in the Southern Ecuadorian Andes: implications for temperature mapping.  
 5209 *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  $6.88 \text{ } ^\circ\text{C}/\text{km}$  for the annual mean temperature.

- 5210  
 5211 [2] Fall, P.L., 1997: Timberline fluctuations and late Quaternary paleoclimates in  
 5212 the Southern Rocky Mountains, Colorado. *Geological Society of America Bul-*  
 5213 *letin*, 109, 1306-1320.

Figure 2a of this paper shows a lapse rate of  $6.9 \text{ } ^\circ\text{C}/\text{km}$  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  $6.0 \text{ } ^\circ\text{C}/\text{km}$  with  $R^2 = 0.80$ .

- 5214  
 5215 [3] Graybill, F.A., H.K. Iyer, 1994: *Regression Analysis: Concepts and Applications*.  
 5216 Duxbury Press, Belmont, California, 701pp.

This book clearly outlines the assumptions of regression. It focuses on concepts and applications about solving practical statistical problems.

- 5217  
 5218 [4] Menne, M.J., C.N. Williams, and R.S. Vose, 2009: The United States Histori-  
 5219 cal Climatology Network monthly temperature data Version 2. *Bulletin of the*  
 5220 *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>

- 5222 [5] Zhang, H.-M., B. Huang, J. Lawrimore, M. Menne, T. M. Smith, 2019:  
 5223 NOAA Global Surface Temperature Dataset (NOAAGlobalTemp), Version  
 5224 5.0 (Time Series). NOAA National Centers for Environmental Information.  
 5225 doi:10.7289/V5FN144H [Access date: March 2021].

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](http://www.4dvd.org).

5226

## Exercises

- 5227 4.1 (a) Compute the temperature lapse rate for August using the TOB mean  
 5228 temperature data from the 24 USHCN stations in Colorado during the period  
 5229 of 1981-2010.  
 5230 (b) Plot the figure similar to Fig. 4.1 for the data and results of Part (a).
- 5231 4.2 Repeat the previous problem but for January. Compare your January results  
 5232 with those of July in Fig. 4.1 and its related text.
- 5233 4.3 (a) Compute the temperature lapse rate for the annual mean temperature  
 5234 based on the 24 USHCN stations in Colorado during the period of 1981-2010.  
 5235 (b) Plot the figure similar to Fig. 4.1 for the data and results of Part (a).
- 5236 4.4 (a) Compute the annual mean temperature lapse rate for a high mountainous  
 5237 region and for a period of your interest.  
 5238 (b) Plot the figure similar to Fig. 4.1 for the data and results of Part (a).
- 5239 4.5 Show that the orthogonality condition Eq. (4.27)

$$\mathbf{e} \cdot \mathbf{x}_a = 0 \quad (4.113)$$

5240 is equivalent to the condition of minimizing SSE given the condition (4.18)

$$\sum_{i=1}^n e_i = 0. \quad (4.114)$$

- 5241 4.6 Show that when  $m = 1$ , the confidence interval formula for the multiple linear  
 5242 regression (4.102) is reduced to the confidence interval formula for the simple  
 5243 linear regression (4.82).  
 4.7 Examine the global average December temperature anomalies from 1880 to  
 5244 2018 in the dataset of the NOAAGlobalTemp.  
 5245 (a) Make a linear regression of the temperature anomalies against time.  
 5246 (b) Compute the confidence intervals of the fitted model at 95% confidence

- 5250 level.
- 5251 (c) Compute the confidence intervals of the anomaly data at 95% confidence
- 5252 level.
- 5253 (d) On the same figure similar to Fig. 4.6(a), plot the scatter plot of the
- 5254 anomaly data against time, and plot the confidence intervals computed in (b)
- 5255 and (c).
- 5256 **4.8** Make a diagnostic analysis for the regression results of the previous problem.
- 5257 (a) Produce a scatter plot of the residuals against time.
- 5258 (b) Visually check whether the assumptions of linearity and constant variance
- 5259 are satisfied.
- 5260 (c) Use the KS-test to check the normality assumption on residuals.
- 5261 (d) Use the DW-test to check the independence assumption on residuals.
- 5262 (e) When serial correlation is considered, find the effective degrees of freedom
- 5263 (edof).
- 5264 (f) Compute the confidence intervals in Steps (b) and (c) in the above
- 5265 problem using edof.
- 5266 (g) Produce a scatter plot of the anomaly data against time, and plot the
- 5267 confidence intervals on the same figure using the results of Step (f).
- 5268 **4.9** (a) Fit the NOAAGlobalTemp's global average annual mean data from 1880
- 5269 to 2018 to a ninth order orthogonal polynomial.
- 5270 (b) Plot the data and the fitted polynomial function on the same figure.
- 5271 (c) Produce a scatter plot of the residuals of the fitting against time.
- 5272 **4.10** Make a diagnostic analysis for the above regression and examine the regres-
- 5273 sion assumptions. In particular, use the KS-test to verify the normality as-
- 5274 sumption, and the DW-test to verify the independence assumption.
- 5275 **4.11** (a) Use multiple linear regression to compute the 12th order polynomial
- 5276 fitting of the NOAAGlobalTemp's global average annual mean data from 1880
- 5277 to 2018:
- $$T = b_0 + b_1 t + b_2 t^2 + \cdots + b_{12} t^{12} + \epsilon. \quad (4.115)$$
- 5278 (b) Plot the data and the fitted polynomial function on the same figure.
- 5279 (c) Produce a scatter plot of the residuals of the fitting against time.
- 5280 **4.12** Make a diagnostic analysis for the above regression and examine the regres-
- 5281 sion assumptions. In particular, use the KS-test to verify the normality as-
- 5282 sumption, and the DW-test to verify the independence assumption.
- 5283 **4.13** (a) Fit the global average **January** monthly mean temperature anomaly
- 5284 data from 1880 to 2018 in the NOAAGlobalTemp dataset to a third order
- 5285 orthogonal polynomial. The global average monthly mean NOAAGlobalTemp
- 5286 time series data are included in the book's master dataset named **data.zip**.
- 5287 You can also download the updated data from the Internet.
- 5288 (b) Plot the data and the fitted polynomial function on the same figure.
- 5289 (c) Produce a scatter plot of the residuals of the fitting against time in another
- 5290 figure.
- 5291 **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.
- 5292
- 5293
- 5294 **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.
- 5295
- 5296
- 5297                   (b) Plot the data and the fitted polynomial function on the same figure.
- 5298                   (c) Produce a scatter plot of the residuals of the fitting against time in another figure.
- 5299
- 5300 **4.16** Make a diagnostic analysis for the previous **July** 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.
- 5301
- 5302
- 5303 **4.17** Use the gridded monthly NOAAGlobalTemp dataset and make a third-order polynomial fit to the **January** monthly mean temperature anomaly data from 1880 to 2018 in a grid box that covers Tokyo, Japan. The gridded monthly NOAAGlobalTemp dataset is included in the book's master data file `data.zip`. You can also download the updated data from the Internet.
- 5304
- 5305
- 5306
- 5307
- 5308 **4.18** Use the gridded monthly NOAAGlobalTemp dataset and make a third-order polynomial fit to the **January** monthly mean temperature anomaly data from 1880 to 2018 in a grid box that covers Bonn, Germany.
- 5309
- 5310
- 5311 **4.19** Use the gridded monthly NOAAGlobalTemp dataset and make a fifth-order polynomial fit to the monthly mean temperature anomaly data for a grid box, a month, and a period of time of your own interest. For example, you may choose your hometown grid box, the month you were born, and the period of 1900-1999.
- 5312
- 5313
- 5314
- 5315
- 5316 **4.20** Make a diagnostic analysis for the previous data fit following the procedures in this chapter.
- 5317
- 5318 **4.21** Use the January time series data of an USHCN station and fit a third order polynomial. Make a diagnostic analysis for the fit. The updated monthly USHCN station data may be downloaded from the Internet.
- 5319
- 5320

5321  
5322 The book uses R and the R Notebook. This chapter explains the installation of R  
5323 and R Studio and demonstrates some basic uses of R.  
5324 Equivalent Python codes and their Jupyter Notebooks may be found at the web-  
5325 site [www.climatemathematics.org](http://www.climatemathematics.org).

## A.1 Download and install R and R-Studio

5329  
5330 For Windows users, visit the website  
5331 <https://cran.r-project.org/bin/windows/base/>  
5332 to find the instructions for R program download and installations.  
5333 For Mac users, visit  
5334 <https://cran.r-project.org/bin/macosx/>  
5335  
5336 If you experience difficulties, please refer to online resources, Google or YouTube.  
5337 A recent 3-minute YouTube instruction for R installation for Windows can be found  
5338 from the following link:  
5339 <https://www.youtube.com/watch?v=0hnk9hcx9M>  
5340  
5341 The same author also has a YouTube instruction about R installation for Mac  
5342 (2 minutes):  
5343 <https://www.youtube.com/watch?v=uXuuWXU-7UQ>  
5344  
5345 When R is installed, one can open R. The R Console window will appear. See  
5346 Fig. A.1. One can use R Console to perform calculations, such as typing 2+3 and  
5347 hitting return. However, most people today prefer using RStudio as the interface.  
5348 To install RStudio, visit  
5349 <https://www.rstudio.com/products/rstudio/download/>  
5350 This site allows one to choose Windows, or Mac OS, or Unix.  
5351  
5352 After both R and RStudio are installed, one can use either R or RStudio, or  
5353 both, depending on one's interest. However, RStudio will not work without R.  
5354 Thus, always install R first.  
5355 When opening RStudio, four windows will appear as shown in Fig. A.2: The  
5356 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.

5357 top of the window can be clicked to run the code. Each run is shown in the lower  
 5358 left R Console window, and recorded on the upper right R History window. When  
 5359 plotting, the figure will appear in the lower right R Plots window. For example,  
 5360 `plot(x,x*x)` renders the eight points in the Plots window, because `x=1:8` defines  
 5361 a sequence of numbers from 1 to 8. `x*x` yields a sequence from  $1^2$  to  $8^2$ .

## A.2 R Tutorial

---

5362 Many excellent tutorials for quickly learning R programming, using a few hours or  
 5363 a few evenings, are available online and in YouTube, such as  
 5364 <http://ww2.coastal.edu/kingw/statistics/R-tutorials/>.  
 5365 You can easily perform an Internet search and find your preferred tutorials.  
 5366 It can be extremely difficult for the beginners of R to navigate through the offi-  
 5367 cial, formal, detailed, and massive R-Project documentation:  
 5368 <https://www.r-project.org/>  
 5369

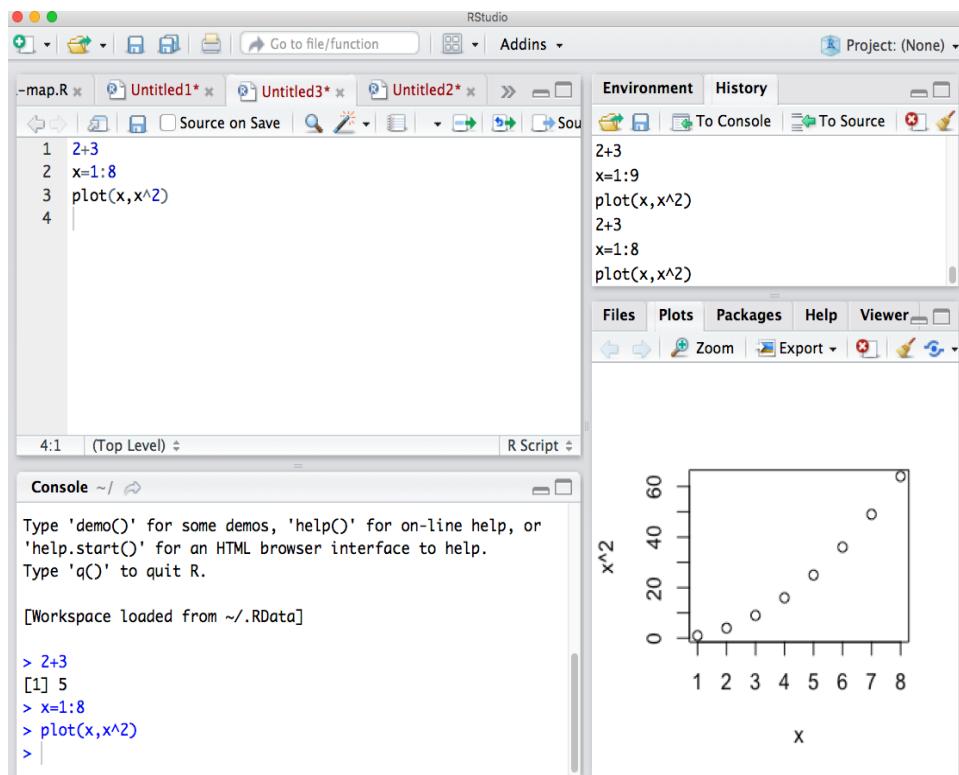


Fig. A.2 R Studio windows.

### A.2.1 R as a smart calculator

R can be used like a smart calculator that allows more elaborate calculations than those done with ordinary calculators.

```

5372 1+4
5373 #[1] 5 This is the result
5374 2+pi/4-0.8
5375 #[1] 1.985398
5376 x<-1
5377 y<-2
5378 z<-4
5379 t<-2*x^y-z
5380 t
5381 #[1] -2
5382 u=2      # "=" sign and "<-" are almost equivalent
5383 v=3      # The text behind the "#" sign is comments
5384 u+v
5385 #[1] 5

```

```
5390 sin(u*v)    # u*v = 6 in the sine function is considered radian by R
5391 #[1] -0.2794155
```

5392 R programming uses assignment operator `a <- b` to assign `b` to `a`. Often the equal  
 5393 operator `a=b` can do the same job or vice versa. The two operators are equivalent  
 5394 in general. However, certain R formulas have specific meanings for `=` and cannot be  
 5395 replaced by `<-`. Most veteran R users use `<-` for assignment and `=` for defined R  
 5396 formulas.

### 5397 A.2.2 Define a sequence in R

---

5399 Directly enter a sequence of daily maximum temperature data at San Diego International  
 5400 Airport (Lat: 32.7336°N, Lon: 117.1831°W) during 1-7 May 2017 [unit:  
 5401 °F].

```
5402 tmax <- c(77, 72, 75, 73, 66, 64, 59)
5403 The data are from the Daily Summary of the Local Climatological Data (LCD),
5404 National Centers for Environmental Information (NCEI)
5405 https://www.ncdc.noaa.gov/cdo-web/datatools/lcd
5406 The command c() is used to hold a data sequence and is named tmax. Entering
5407 the tmax command will render the temperature data sequence:
```

```
5408 tmax
5409 #[1] 77 72 75 73 66 64 59
```

5410 You can generate different sequences using R, e.g.,
 5411 `1: 8` #Generates a sequence 1,2,...,8

5412 Here the pound sign `#` begins R comments which are not executed by R calculations.  
 5413 The same sequence can be generated by different commands, such as

```
5414 seq(1,8)
5415 seq(8)
5416 seq(1,8, by=1)
5417 seq(1,8, length=8)
5418 seq(1,8, length.out =8)
```

5419 The most useful sequence commands are `seq(1,8, by=1)` and `seq(1,8, length=8)`  
 5420 or `seq(1,8, len=8)`. The former is determined by a begin value, end value, and  
 5421 step size, and the latter by a begin value, end value, and number of values in the  
 5422 sequence. For example, `seq(1951,2016, len=66*12)` renders a sequence of all the  
 5423 months from January 1951 to December 2016.

---

### A.2.3 Define a function in R

---

5424  
5425  
5426 The function command is  
5427 `name <- function(var1, var2, ...) expression of the function.`  
5428 For example,

```
5429 samfctn <- function(x) x*x
5430 samfctn(4)
5431 #[1] 16
5432 fctn2 <- function(x,y,z) x+y-z/2
5433 fctn2(1,2,3)
5434 #[1] 1.5
```

---

### A.2.4 Plot with R

---

5435  
5436  
5437 R can plot all kinds of curves, surfaces, statistical plots, and maps. Below are a  
5438 few very simple examples for R beginners. For adding labels, ticks, color, and other  
5439 features to a plot, you will learn them from later parts of this book, and you may  
5440 also carry out an Internet search on R plot to find the commands for the proper  
5441 inclusion of the desired features.

5442 R plotting is based on the coordinate data. The following command plots the  
5443 seven days of San Diego Tmax data above:

```
5444 plot(1:7, c(77, 72, 75, 73, 66, 64, 59))
```

The resulting graph is shown in Fig. A.3.

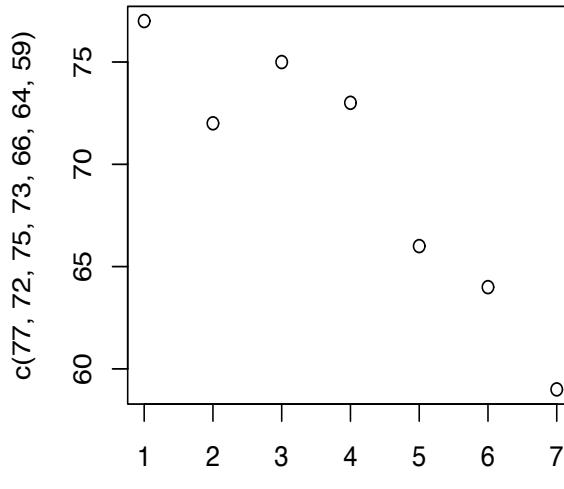


Fig. A.3

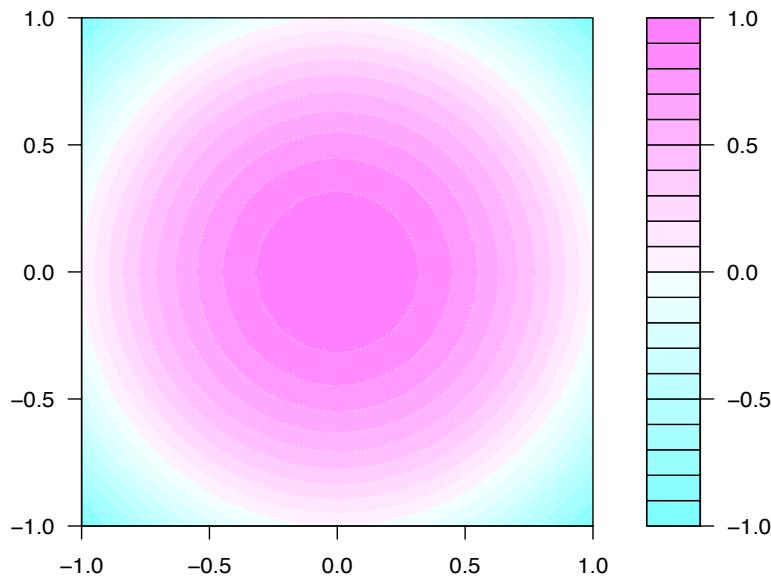
The daily maximum temperature during 1-7 May 2017 at San Diego International Airport.

```

5446 plot(sin, -pi, 2*pi)    #plot the curve of y=sin(x) from -pi to 2 pi
5447
5448 square <- function(x) x*x  #Define a function
5449 plot(square, -3,2)      # Plot the defined function
5450
5451 # Plot a 3D surface
5452 x <- seq(-1, 1, length=100)
5453 y <- seq(-1, 1, length=100)
5454 z <- outer(x, y, function(x, y)(1-x^2-y^2))
5455 #outer (x,y, function) renders z function on the x, y grid
5456 persp(x,y,z, theta=330)
5457 # yields a 3D surface with perspective angle 330 deg
5458
5459 #Contour plot\index{contour plot}
5460 contour(x,y,z) #lined contours
5461 filled.contour(x,y,z) #color map of contours

```

5462 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  $z = 1 - x^2 - y^2$ .

5463

## A.2.5 Symbolic calculations by R

---

5464 Many people once thought that R can handle only numbers. Actually R can also do  
 5465 symbolic calculations, such as finding a derivative, although at present R is not the

best symbolic calculation tool. One can use WolframAlpha, SymPy, and Yacas for free symbolic calculations or use the paid software packages Maple or Mathematica. Carry out an Internet search on symbolic calculation for calculus to find a long list of symbolic calculation software packages, e.g.,  
[https://en.wikipedia.org/wiki/List\\_of\\_computer\\_algebra\\_systems](https://en.wikipedia.org/wiki/List_of_computer_algebra_systems).

```

5468 D(expression(x^2,'x'), 'x')
5469 # Take derivative of x^2 w.r.t. x
5470 2 * x #The answer is 2x
5471
5472
5473 fx= expression(x^2,'x') #assign a function
5474 D(fx,'x') #differentiate the function w.r.t. x
5475 2 * x #The answer is 2x
5476
5477
5478 fx= expression(x^2*sin(x),'x')
5479 #Change the expression and use the same derivative command
5480 D(fx,'x')
5481 2 * x * sin(x) + x^2 * cos(x)
5482
5483
5484 fxy = expression(x^2+y^2, 'x','y')
5485 #One can define a function of 2 or more variables
5486 fxy #renders an expression of the function in terms of x and y
5487 #expression(x^2 + y^2, "x", "y")
5488 D(fxy,'x') #yields the partial derivative with respect to x: 2 * x
5489 D(fxy,'y') #yields the partial derivative with respect to y: 2 * y
5490
5491
5492 square = function(x) x^2
5493 integrate (square, 0,1)
5494 #Integrate x^2 from 0 to 1 equals to 1/3 with details below
5495 #0.3333333 with absolute error < 3.7e-15
5496
5497
5498 integrate(cos,0,pi/2)
5499 #Integrate cos(x) from 0 to pi/2 equals to 1 with details below
5500 #1 with absolute error < 1.1e-14

```

5501 The above two integration examples are for definite integrals. It seems that no  
5502 efficient R packages are available for finding antiderivatives, or indefinite integrals.

## A.2.6 Vectors and matrices

5503  
5504 R can handle all kinds of operations involving vectors and matrices.

```

5505 c(1,6,3,pi,-3) #c() gives a vector and is considered a 4X1 column vector\index{column vec}
5506 #[1] 1.000000 6.000000 3.000000 3.141593 -3.000000

```

```
5508 seq(2,6) #Generate a sequence from 2 to 6
5509 #[1] 2 3 4 5 6
5510 seq(1,10,2) # Generate a sequence from 1 to 10 with 2 increment
5511 #[1] 1 3 5 7 9
5512 x=c(1,-1,1,-1)
5513 x+1 #1 is added to each element of x
5514 #[1] 2 0 2 0
5515 2*x #2 multiplies each element of x
5516 #[1] 2 -2 2 -2
5517 x/2 # Each element of x is divided by 2
5518 #[1] 0.5 -0.5 0.5 -0.5
5519 y=seq(1,4)
5520 x*y # This multiplication * multiples each pair of elements
5521 #[1] 1 -2 3 -4
5522 x%*%y #This is the dot product of two vectors and yields
5523 # [,1]
5524 #[1,] -2
5525 t(x) # Transforms x into a row 1X4 vector
5526 # [,1] [,2] [,3] [,4]
5527 #[1,] 1 -1 1 -1
5528 t(x)%*%y #This is equivalent to dot product and forms 1X1 matrix
5529 # [,1]
5530 #[1,] -2
5531 x%*%t(y) #This column times row yields a 4X4 matrix\index{matrix}
5532 # [,1] [,2] [,3] [,4]
5533 #[1,] 1 2 3 4
5534 #[2,] -1 -2 -3 -4
5535 #[3,] 1 2 3 4
5536 #[4,] -1 -2 -3 -4
5537 my=matrix(y,ncol=2)
5538 #Convert a vector into a matrix of the same number of elements
5539 #The matrix elements\index{matrix!elements} go by column, the first column, second, etc
5540 #Command matrix(y,ncol=2, nrow=2), or matrix(y, ncol=2)
5541 #or matrix(y,2), or matrix(y,2,2) does the same job
5542 my
5543 # [,1] [,2]
5544 #[1,] 1 3
5545 #[2,] 2 4
5546 dim(my) #find dimensions of a matrix
5547 #[1] 2 2
5548
5549 bigM=matrix(1:100, nrow=10) #Generate a 10-by-10 matrix
5550 subM=bigM[4:6,3:7] #Extract a sub-matrix from a big matrix
5551 subM
```

```
5552 #      [,1] [,2] [,3] [,4] [,5]
5553 #[1,] 24   34   44   54   64
5554 #[2,] 25   35   45   55   65
5555 #[3,] 26   36   46   56   66
5556
5557 as.vector(my) #Convert a matrix to a vector, again via columns
5558 #[1] 1 2 3 4
5559 mx <- matrix(c(1,1,-1,-1), byrow=TRUE,nrow=2)
5560 mx*my #multiplication between each pair of elements
5561 #      [,1] [,2]
5562 #[1,] 1    3
5563 #[2,] -2   -4
5564 mx/my #division between each pair of elements
5565 #      [,1]      [,2]
5566 #[1,] 1.0  0.3333333
5567 #[2,] -0.5 -0.2500000
5568 mx-2*my
5569 #      [,1] [,2]
5570 #[1,] -1   -5
5571 #[2,] -5   -9
5572 mx%*%my #This is the real matrix multiplication\index{matrix multiplication} in matrix th
5573 #      [,1] [,2]
5574 #[1,] 3    7
5575 #[2,] -3   -7
5576 det(my) #determinant\index{determinant}
5577 #[1] -2
5578 myinv = solve(my) #yields the inverse of a matrix
5579 myinv
5580 #      [,1] [,2]
5581 #[1,] -2   1.5
5582 #[2,] 1    -0.5
5583 myinv%*%my #verifies the inverse of a matrix
5584 #      [,1] [,2]
5585 #[1,] 1    0
5586 #[2,] 0    1
5587 diag(my) #yields the diagonal vector of a matrix
5588 #[1] 1 4
5589 myeig=eigen(my) #yields eigenvalues\index{eigenvalue} and unit eigenvectors\index{eigenve
5590 myeig
5591 myeig$values
5592 #[1] 5.3722813 -0.3722813
5593 myeig$vectors
5594 #      [,1]      [,2]
5595 #[1,] -0.5657675 -0.9093767
```

```

5596 #[2,] -0.8245648 0.4159736
5597 mysvd = svd(my) #SVD decomposition of a matrix M=UDV'
5598 #SVD can be done for a rectangular matrix of mXn
5599 mysvd$d
5600 #[1] 5.4649857 0.3659662
5601 mysvd$u
5602 # [,1] [,2]
5603 #[1,] -0.5760484 -0.8174156
5604 #[2,] -0.8174156 0.5760484
5605 mysvd$v
5606 # [,1] [,2]
5607 #[1,] -0.4045536 0.9145143
5608 #[2,] -0.9145143 -0.4045536
5609
5610 ysol=solve(my,c(1,3))
5611 #solve linear equations\index{linear equations} matrix %*% x = b
5612 ysol #solve(matrix, b)
5613 #[1] 2.5 -0.5
5614 my%*%ysol #verifies the solution
5615 # [,1]
5616 #[1,] 1
5617 #[2,] 3

```

---

### A.2.7 Simple statistics by R

---

5620 R was originally designed to do statistical calculations. Thus, R has a comprehensive  
 5621 set of statistics functions and software packages. This sub-section gives a few basic  
 5622 commands. More will be described in the statistics chapter of this book.

```

5623 x=rnorm(10) #generate 10 normally distributed numbers
5624 x
5625 #[1] 2.8322260 -1.2187118 0.4690320 -0.2112469 0.1870511
5626 #[6] 0.2275427 -1.2619005 0.2855896 1.7492474 -0.1640900
5627 mean(x)
5628 #[1] 0.289474
5629 var(x)
5630 #[1] 1.531215
5631 sd(x)
5632 #[1] 1.237423
5633 median(x)
5634 #[1] 0.2072969
5635 quantile(x)
5636 # 0% 25% 50% 75% 100%
5637 #-1.2619005 -0.1994577 0.2072969 0.4231714 2.8322260

```

```

5638 range(x) #yields the min and max of x
5639 #[1] -1.261900 2.832226
5640 max(x)
5641 #[1] 2.832226
5642
5643 boxplot\index{boxplot}(x) #yields the box plot of x
5644 w=rnorm(1000)
5645
5646 summary(rnorm(12)) #statistical summary of the data sequence
5647 # Min. 1st Qu. Median Mean 3rd Qu. Max.
5648 #-1.9250 -0.6068 0.3366 0.2309 1.1840 2.5750
5649
5650 hist(w)
5651 #yields the histogram\index{histogram} of 1000 random numbers with a normal distribution
5652
5653 #Linear regression\index{linear regression} and linear trend\index{linear trend} line
5654 #2007–2016 data of the global temperature\index{anomalies!temperature} anomalies
5655 #Source: NOAAGlobalTemp data
5656 t=2007:2016
5657 T=c(.36,.30, .39, .46, .33, .38, .42, .50, .66, .70)
5658 lm(T ~ t) #Linear regression model of temp vs time
5659 #(Intercept) t
5660 #-73.42691 0.03673
5661 #Temperature change rate is 0.03673 deg C/yr or 0.37 deg C/decade
5662 plot(t,T, type="o", xlab="Year", ylab="Temperature [deg C]",
5663 main="2007–2016 Global Temperature\index{anomalies!temperature} Anomalies\index{anom
5664 and Their Linear Trend [0.37 deg C/decade] ")
5665 abline(lm(T ~ t), lwd=2, col="red") #Regression line

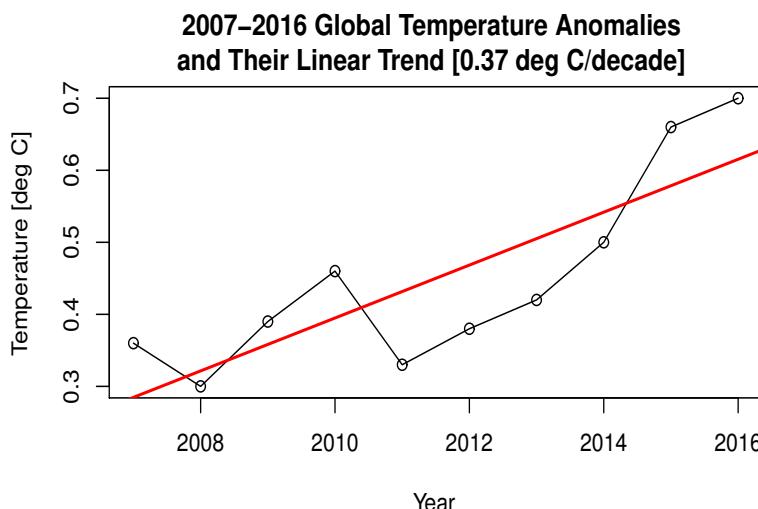
```

5666 The global temperature data from 2007–2016 in the above R code example are  
 5667 displayed in Fig. A.5, together with their linear trend line

$$T = -73.42691 + 0.03673t. \quad (\text{A.1})$$

### 5668 A.3 Online Tutorials

5670 Numerous online R tutorials are available. Several are relatively efficient for learning  
 5671 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.

5672  
5673

### A.3.1 YouTube tutorial: for true beginners

---

5674 This is a very good and slow-paced 22-minute YouTube tutorial: Chapter 1. An  
 5675 Introduction to R  
 5676 <https://www.youtube.com/watch?v=suVFuGET-OU>

5677  
5678

### A.3.2 YouTube tutorial: for some basic statistical summaries

---

5679 This is a 9-minute tutorial by Layth Alwan.  
 5680 <https://www.youtube.com/watch?v=XjOZQN-Nre4>

5681  
5682

### A.3.3 YouTube tutorial: Input data by reading a csv file into R

---

5683 An excel file can be saved as csv file: xxxx.csv. This 15-minute YouTube video by  
 5684 Layth Alwan shows how to read a csv file into R. He also shows linear regression.  
 5685 <https://www.youtube.com/watch?v=QkE8cp0B9gg>

5686

5687 R can input all kinds of data files, including xlsx, txt, netCDF, MatLab data,  
 5688 Fortran file, and SAS data. Some commands are below. One can search the Internet  
 5689 to find proper data reading commands for any particular data format.

```
5690 mydata <- read.csv("mydata.csv")
5691 # read csv file named "mydata.csv"
5692
```

```

5693 mydata <- read.table("mydata.txt")
5694 # read text file named "my data.txt"
5695
5696 library(gdata)           # load the gdata package
5697 mydata = read.xls("mydata.xls") # read an excel file
5698
5699 library(foreign)          # load the foreign package
5700 mydata = read.mtp("mydata.mtp") # read from .mtp file
5701
5702 library(foreign)          # load the foreign package
5703 mydata = read.spss("myfile", to.data.frame=TRUE)
5704
5705 ff <- tempfile()
5706 cat(file = ff, "123456", "987654", sep = "\n")
5707 read.fortran(ff, c("F2.1","F2.0","I2")) #read a fotran file
5708
5709 library(ncdf4)
5710 ncin <- ncdf4::nc_open("ncfname") # open a NetCDF file
5711 lon <- ncvar_get(ncin, "lon") #read data "lon" from a netCDF file into R

```

5712 Many more details of reading and reformatting of .nc files will be discussed later  
 5713 when dealing with NCEP/NCAR Reanalysis data.

5714 Some libraries are not in the R project. For example,

```

5715 library(ncdf4) #The following error message pops up
5716 Error in library(ncdf4) : there is no package called ncdf4

```

5717 You can install the R package by

```

5718 install.packages("ncdf4")

```

5719 After this installation, `library(ncdf4)` will run, and the functions in the ncdf4  
 5720 package will work.

5721 You only need to install the package once on your computer, but in each new R  
 5722 session you must run `library(package)` in order to activate the package functions.  
 5723 Many examples will be shown in the rest of this book.

5724 The R packages and the datasets used in this book are listed below and can be  
 5725 downloaded and installed first before proceeding to the R codes in the rest of the  
 5726 book.

```

5727 #R packages: animation, chron, e1071, fields, ggplot2, lattice,
5728 #latticeExtra, maps, mapdata, mapproj, matrixStats, ncdf4,
5729 #NLRoot, RColorBrewer, rgdal, rasterVis, raster, sp, TTR
5730
5731 #The zipped data:

```

```

5732 #https://cambridge.org/climate mathematics/data.zip
5733
5734 #To load a single package, such as "animation", you can do
5735 library(animation)
5736
5737 #You can also load all these packages in one shot using
5738 # pacman
5739
5740 install.packages("pacman")
5741 library(pacman)
5742 pacman::p_load(animation, chron, e1071, fields, ggplot2, lattice,
5743                  latticeExtra, maps, mapdata, mapproj, matrixStats, ncdf4,
5744                  NLRoot, RColorBrewer, rgdal, rasterVis, raster, sp, TTR)

```

5745 On your computer, you can create a directory called `climmath` under your user  
 5746 name. The one used in the book is `Users/sshenn/climmath`. You unzip the data and  
 5747 move the data folder under the `Users/sshenn/climmath` directory. A data folder will  
 5748 created: `Users/sshenn/climmath/data`. The data folder contains about 400 MB of  
 5749 data. Place all the R codes in the directory `Users/sshenn/climmath`. Then, you can  
 5750 run all the codes in this book after replacing `sshenn` by your user name on your own  
 5751 computer.

## 5752 A.4 Chapter summary

---

5753 This chapter has described the following basic aspects of the R programming lan-  
 5754 guage

- 5755 (i) Installation of both R and R Studio.
- 5756 (ii) Layout of R Studio panels and functions.
- 5757 (iii) R commands for defining and computing vectors, matrices, and functions.
- 5758 (iv) Simple statistics using R: mean, median, standard deviation, variance, his-  
5759 togram, boxplot, scatterplot, and linear regression.
- 5760 (v) Online resources for R tutorial material.

5761 We suggest that reading and practicing the R basics included this chapter might  
 5762 require about three hours. You might then use perhaps five to ten hours to do  
 5763 the exercise problems of this chapter. After that introduction to R, we estimate  
 5764 that you should be able to develop and carry out projects using R independently.  
 5765 Familiarizing yourself with additional R commands may require spending some time  
 5766 online with a search engine.

5767 More sophisticated examples of using R for the analysis and visualization of the  
 5768 data from climate models and observations are included in Chapters 9-11. These

5770 chapters also provide research-level R-graphics codes for climate sciences, R pro-  
5771 grams for advanced statistical analysis of climate data, such as empirical orthogonal  
5772 functions computed from the climate model datasets, and R programming tech-  
5773 niques for analyzing datasets with missing data, either in space or time. To develop  
5774 R projects at a more advanced level, you will need to read these chapters or search  
5775 for the R code for the specific analysis or graphics tasks described in these chapters.  
5776 These tasks include examples such as the preparation of data from a global climate  
5777 model (GCM) for a singular value decomposition (SVD) analysis.

5778 The R programming language was created by statisticians Ross Ihaka and Robert  
5779 Gentleman in New Zealand and was first released in 1995. R is named partly after  
5780 the first names of the two original two authors of R, and partly as a play on the  
5781 name of the S programming language for statistics. Currently, R is a popular com-  
5782 puter programming language used in almost every field of science and engineering.  
5783 Among the top programming languages ranked by the Institute of Electric and  
5784 Electronic Engineers (IEEE) in 2017, R ranks sixth, following Python, C, Java,  
5785 C++, and C#. For a climate science student or professional who is not a specialist  
5786 in computer programming or information technology, R is easy to learn, and it offers  
5787 numerous public-domain software packages that are free to use. As an alternative  
5788 to R, equivalent Python codes for our entire book are available online at the book  
5789 website [www.climatemathematics.org](http://www.climatemathematics.org).

## References and Further Readings

5791 [1] King, K.B., 2016: *R Tutorials*, Coastal Carolina University. URL:

5792 <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.

5794 [2] Jost, S., 2018: *Introduction to R: An R Tutorial for Data Analysis and Regres-*  
5795 *sion*. De Paul University. URL:

5796 <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.

5798 [3] Vallis, G.K., 2016: Geophysical fluid dynamics: whence, whither  
5799 and why? *Proceedings of the Royal Society A*, 472: 20160140.  
5800 <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.

5802

## Exercises

5803

5804   **A.1** Use R to define a data sequence `t=seq(2015,2018, length=100)`, and then  
 5805    plot the following two functions on the same figure:  $y = \sin(2\pi(t - 0.1))$  and  
 5806     $y = \cos^2(2\pi t)$ .

5807   **A.2** (a) Use R to make a contour plot of the function  $z = \sin^2 x \cos^2(y - \pi)$  over  
 5808    the domain of  $[0, 2\pi] \times [0, 2\pi]$ .  
 5809    (b) Use R to plot a color contour map for the same function on the same  
 5810    domain.

5811   **A.3** Use R to solve the following linear equations:

$$\begin{cases} 9x + 8y = 87 \\ 6x - 20y = 126 \end{cases}$$

5812   **A.4** Use R to solve the following linear equations:

$$\begin{cases} -3x + 2y + z = 1 \\ -2x - y + z = 2 \\ 2x + y - 4z = 0 \end{cases}$$

5813   **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.

5832    Go to the USHCN website

5833    [http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn\\_map\\_interface.html](http://cdiac.ornl.gov/epubs/ndp/ushcn/ushcn_map_interface.html)  
 5834

5835    and download the monthly Tmax, Tmin, and Tmean data of the Cuyamaca

5836 station (USHCN Site No. 042239) near San Diego, California, USA. Or down-  
5837 load the data of this book at

5838 [www.cambridge.org/climate-mathematics/data.zip](http://www.cambridge.org/climate-mathematics/data.zip)

5839

5840 Unzip the data and find the data from the file named CA042239T.csv

5841 (a) Use R to arrange the monthly Cuyamaca Tmax data from January 1961  
5842 to December 1990 as a matrix with each row as year and each column as  
5843 month.

5844 (b) Do the same for Tmin.  
5845 (c) Do the same for Tmean.

5846 **A.6** (a) Use R to calculate the August climatology of Tmax, Tmin, and Tmean  
5847 for the Cuyamaca station according to the 1961-1990 climatology period.

5848 (b) Use R to compute the standard deviation of Tmax, Tmin, and Tmean of  
5849 the Cuyamaca station for January during the 1961-1990 climatology period.

5850 **A.7** (a) Use R to plot the the Cuyamaca January Tmin time series from 1951 to  
5851 2010 with a continuous curve.

5852 (b) Use R to plot the linear trend lines of Tmin on the same plot as (a) in  
5853 the following time periods:

5854 (i) 1951-2010,  
5855 (ii) 1961-2010,  
5856 (iii) 1971-2010, and  
5857 (iv) 1981-2010.

5858 (c) Finally, what is the temporal trend per decade for each of the four periods  
5859 above?

5860 **A.8** Use R to plot the time series and its trend line for P.D. Jones' global average  
5861 annual mean temperature anomaly data: JonesGlobalT.txt . This data file  
5862 can be found from the book's data.zip file downloaded from the book website.

5863 (a) Plot the global average annual mean temperature from 1880 to 2015.

5864 (b) Find the linear trend of the temperature from 1880 to 2015. Plot the  
5865 trend line on the same figure as (a).

5866 (c) Find the linear trend from 1900 to 1999. Plot the trend line on the same  
5867 figure as (a).

5868 **A.9** Use the gridded NOAA global monthly temperature anomaly data NOAA-  
5869 GlobalTemp from the following website or another data source

5870 [https://www.ncdc.noaa.gov/data-access/marineocean-data/  
5871 noaa-global-surface-temperature-noaaglobaltemp](https://www.ncdc.noaa.gov/data-access/marineocean-data/noaa-global-surface-temperature-noaaglobaltemp)

5872 Or use the NOAAGlobalT.csv data file from the book's data.zip file down-  
5873 loaded from the book website. Choose two 5-by-5 degrees lat-lon grid boxes  
5874 of your interest. Plot the temperature anomaly time series of the two boxes  
5875 on the same figure using two different colors.

- 5876   **A.10** Use the same NOAAGlobalTemp dataset, choose sufficiently many grid boxes  
5877    that cover the state of Texas, USA. Compute the average temperature anomalies  
5878    of these boxes for each month. Then plot the monthly average temperature  
5879    anomalies as a function of time. Plot a linear trend line on the same figure.  
5880   **A.11** Choose a 5-by-5 degrees grid box that covers Edmonton, Canada, and another  
5881    grid box that covers San Diego, USA.  
5882      (a) Use R and 30 years of the January NOAAGlobalTemp data from January  
5883       1981 to January 2010 to compute the standard deviations for each grid box  
5884       for January.  
5885      (b) Do the same for February, March, ..., December.  
5886      (c) Use R to write your standard deviation results in a 12-by-2 matrix with  
5887       each row for a month, and each column for a grid box ID.  
5888      (d) Describe the main differences between the values of the two columns.

5889

**B**

## Appendix B Visualization of Matrices

5891

5892 People talk about climate data frequently, also read or imagine climate data, and  
 5893 yet rarely play with them and use them, because people often think that it takes  
 5894 a computer expert to do that. However, that has changed. With today's technol-  
 5895 ogy, now anyone can use a computer to play with climate data, such as a sequence  
 5896 of temperature values of a weather station at different observed times, a matrix  
 5897 of data for a station for temperature, air pressure, precipitation; wind speed, and  
 5898 wind direction at different times; and an array of temperature data on a 5-degree  
 5899 latitude-longitude grid for the entire world for different months. The first is a vec-  
 5900 tor. The second is a variable-time matrix, and a space-time 3-dimensional array.  
 5901 When considering temperature variation in time at different air pressure levels and  
 5902 different water depth, we need to add one more dimension: the altitude. The tem-  
 5903 perature data for ocean and atmosphere for the Earth is a 4-dimensional array,  
 5904 with 3D space and 1D time. This chapter attempts to provide basic statistical and  
 5905 computing methods to describe and visualize some simple climate datasets. As the  
 5906 book progresses, more complex statistics and data visualization will be introduced.

5907 We use both R and Python computer codes in this book for computing and  
 5908 visualization. Our method description is stated in R. A Python code following each  
 5909 R code is included in a box with a light yellow background. You can also learn  
 5910 the two computer languages and their applications to climate data from the book  
 5911 “*Climate Mathematics: Theory and Applications*” (Shen and Somerville 2019) and  
 5912 its website [www.climatemathematics.org](http://www.climatemathematics.org).

5913 The climate data used in this book are included in the `data.zip` file download-  
 5914 able from our book website [www.climatestatistics.org](http://www.climatestatistics.org). You can also obtain the  
 5915 updated data from the original data providers, such as [www.esrl.noaa.gov](http://www.esrl.noaa.gov) and  
 5916 [www.ncei.noaa.gov](http://www.ncei.noaa.gov).

5917 After learning this chapter, a reader should be able to analyze simple climate  
 5918 datasets, compute data statistics, and plot the data in various ways.

5919

### B.1 Global temperature anomalies from 1880 to 2018: data visualization and statistical indices

5920

5922 In a list of popular climate datasets, the global average annual mean surface air  
 5923 temperature anomalies might be on top. Here, the *anomalies* means the temperature

5924 departures from the normal temperature that is called *climatology*. Climatology is  
 5925 usually defined as the mean of temperature data in a given period of time, such as  
 5926 from 1971 to 2020. Thus, the temperature anomaly data are the differences of the  
 5927 temperature data minus the climatology.

5928 This section will use the global average annual mean surface air temperature  
 5929 anomaly dataset as an example to describe some basic statistical and computing  
 5930 methods.

### 5931 **B.1.1 The NOAAGlobalTemp dataset**

---

5933 The 1880-2018 global average annual mean surface air temperature (SAT) anomaly  
 5934 data are shown as follows:

```
5935 [1] -0.37 -0.32 -0.32 -0.40 -0.46 -0.47 -0.45 -0.50 -0.40 -0.35 -0.57
5936 [12] -0.49 -0.55 -0.56 -0.53 -0.47 -0.34 -0.36 -0.50 -0.37 -0.31 -0.39
5937 [23] -0.49 -0.58 -0.66 -0.53 -0.46 -0.62 -0.69 -0.68 -0.63 -0.68 -0.58
5938 [34] -0.57 -0.39 -0.32 -0.54 -0.56 -0.45 -0.45 -0.46 -0.39 -0.47 -0.46
5939 [45] -0.50 -0.39 -0.31 -0.40 -0.42 -0.54 -0.34 -0.32 -0.36 -0.49 -0.35
5940 [56] -0.38 -0.36 -0.26 -0.27 -0.26 -0.15 -0.05 -0.09 -0.09 0.04 -0.08
5941 [67] -0.25 -0.30 -0.30 -0.30 -0.41 -0.26 -0.22 -0.15 -0.36 -0.38 -0.44
5942 [78] -0.19 -0.13 -0.18 -0.22 -0.16 -0.15 -0.13 -0.39 -0.32 -0.27 -0.26
5943 [89] -0.27 -0.15 -0.21 -0.32 -0.22 -0.08 -0.32 -0.24 -0.32 -0.05 -0.13
5944 [100] -0.02 0.02 0.06 -0.06 0.10 -0.10 -0.11 -0.01 0.13 0.13 0.05
5945 [111] 0.19 0.16 0.01 0.03 0.09 0.21 0.08 0.27 0.39 0.20 0.18
5946 [122] 0.30 0.35 0.36 0.33 0.41 0.37 0.36 0.30 0.39 0.45 0.33
5947 [133] 0.38 0.42 0.50 0.66 0.70 0.60 0.54
```

5948 The data are part of the dataset named as the NOAA Merged Land Ocean Global  
 5949 Surface Temperature Analysis (NOAAGlobalTemp) V4. The dataset was generated  
 5950 by the NOAA National Centers for Environmental Information (NCEI) (Smith et  
 5951 al. 2008; Vose et al. 2012). Here, NOAA stands for the United States National  
 5952 Oceanic and Atmospheric Administration.

5953 The anomalies are with respect to the 1971-2000 climatology, i.e., 1971-2000  
 5954 mean. An anomaly of a weather station datum is defined by the datum minus the  
 5955 station climatology. The first anomaly datum  $-0.37^{\circ}\text{C}$ , indexed by [1], in the  
 5956 above data table corresponds to 1880 and the last to 2018, a total of 139 years. The  
 5957 last row is indexed from [133] to [139].

5958 One might be interested in various kinds of statistical characteristics of the data,  
 5959 such as mean, variance, standard deviation, skewness, kurtosis, median, 5th per-  
 5960 centile, 95th percentile, and other quantiles. Is the data's probabilistic distribution  
 5961 approximately normal? What does the box plot look like? Are there any outliers?  
 5962 What is a good graphic representation of the data, i.e., popularly known as a climate  
 5963 figure?

5964 When considering global climate changes, why do scientists often use anomalies,

5965 instead of the full values directly from the observed thermometer readings? This is  
 5966 because that the observational estimates of the global average annual mean surface  
 5967 temperature are less accurate than the similar estimates for the changes from year  
 5968 to year. There is a concept of characteristic spatial correlation length scale for a  
 5969 climate variable, such as surface temperature. The length scale is often computed  
 5970 from anomalies.

5971 The use of anomalies is also a way of reducing or eliminating individual station  
 5972 biases. A simple example of such biases is that due to station location, which is  
 5973 usually fixed in a long period of time. It is easy to understand, for instance, that  
 5974 a station located in the valley of a mountainous region might report surface tem-  
 5975 peratures that are higher than the true average surface temperature for the entire  
 5976 region and cannot be used to describe the behavior of climate change in the region.  
 5977 However, the anomalies at the valley station may synchronously reflect the charac-  
 5978 teristics of the anomalies for the region. Many online materials give justifications  
 5979 and examples on the use of climate data anomalies, e.g., NOAA NCEI (2021).

5980 The global average annual mean temperature anomalies quoted above are also  
 5981 important for analyzing climate simulations. When we average over a large scale,  
 5982 many random errors cancel out. When we investigate the response of such large  
 5983 scale perturbations as the variations of Sun's brightness or atmospheric carbon  
 5984 dioxide, these averaged data can help validate and improve climate models. See the  
 5985 examples in the book by Hennemuth et al. (2013) that includes many statistical  
 5986 analyses of both observed and model data.

5987 **B.1.2 Visualize the data of global average annual mean**  
 5988 **temperature**

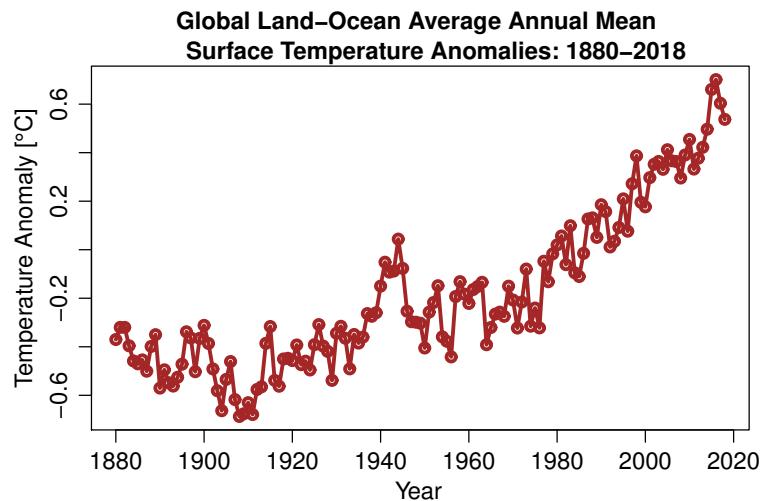
---

5990 Many different ways have been employed to visualize the global average annual  
 5991 mean temperature anomalies. The following three are popular ones appearing in  
 5992 scientific and news publications: (a) a simple point-line graph, (b) a curve of stair-  
 5993 case steps, and (c) a color bar chart, as shown in Figs. B.1 - B.3. This subsection  
 5994 shows how to generate these figures by R and Python computer programming lan-  
 5995 guages. The Python codes are in yellow boxes.

5996 **B.1.2.1 Plot a point-line graph of time series data**

5997 Figure B.1 is a simple line graph that connects all the data points, denoted by  
 5998 the small circles, together to form a curve showing the historical record of the  
 5999 global temperature anomalies. It is plotted from the the NOAAGlobalTemp V4  
 6000 data quoted above. The figure can be generated by the following R code.

```
6001 # R plot Fig. 1.1: A simple line graph of data
6002 # go to your working directory
6003 setwd("/Users/sshen/climstats")
6004 # read the data file from the folder named "data"
6005 NOAAtemp = read.table(
6006   "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.

```

6007   header=FALSE) #Read from the data folder
6008   dim(NOAAtemp) # check the data matrix dimension
6009   #[1] 140 6 #140 years from 1880 to 2019
6010   #2019 will be excluded since data only up to July 2019
6011   #col1 is year, col2 is anomalies, col3-6 are data errors
6012   par(mar=c(3.5,3.5,2.5,1), mgp=c(2,0.8,0))
6013   plot(NOAAtemp[1:139,1], NOAAtemp[1:139,2],
6014     type = "l", col="brown", lwd=3, cex.lab=1.2, cex.axis=1.2,
6015     main ="Global\u2014Land\u2014Ocean\u2014Average\u2014Annual\u2014Mean
6016     Surface\u2014Temperature\u2014Anomalies\u20141880\u20142018",
6017     xlab="Year",
6018     ylab=expression(
6019       paste("Temperature\u2014Anomaly\u2014[", degree, "C]")))

```

```

# Python plot Fig. 1.1: A simple line graph of data
# Go to your working directory
os.chdir("/Users/sshen/climstats")
# read the data file from the folder named "data"
NOAAtemp = read_table(
    "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

```

6020

6021

### B.1.2.2 Plot a staircase chart

6022 The staircase chart Fig. B.2 shows both data and the year-to-year annual changes  
 6023 of temperature. One can clearly see that some years have a larger change from  
 6024 their previous year, while other years have a smaller change. This information can  
 6025 be used for further climate analysis.

6026 Denote the global average annual mean temperature by  $T(t)$  where  $t$  stands for  
 6027 time in year, the 1971-2000 climatology is by  $C$ , and the anomalies by  $A(t)$ . We  
 6028 have

$$T(t) = C + A(t). \quad (\text{B.1})$$

6029 The temperature change from its previous year is

$$\Delta T = T(t) - T(t-1) = A(t) - A(t-1) = \Delta A. \quad (\text{B.2})$$

6030 This implies that an anomaly change is equal to its corresponding temperature  
 6031 change. We do not need to evaluate the 1971-2000 climatology  $C$  of the global  
 6032 average annual mean, when studying changes, as discussed earlier.

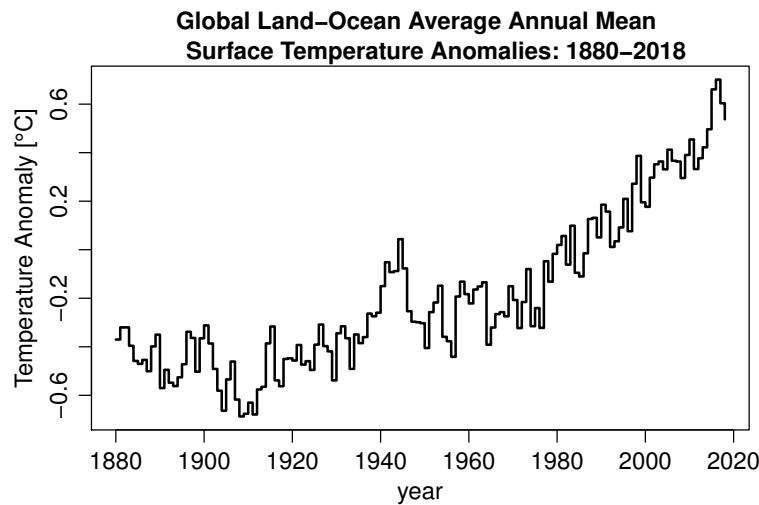
6033 Figure B.2 can be generated by the following R code

```

#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]")))

```

6044 )



**Fig. B.2** Staircase chart of the 1880–2018 global average annual mean temperature anomalies.

6045 The corresponding Python code is in the following box.

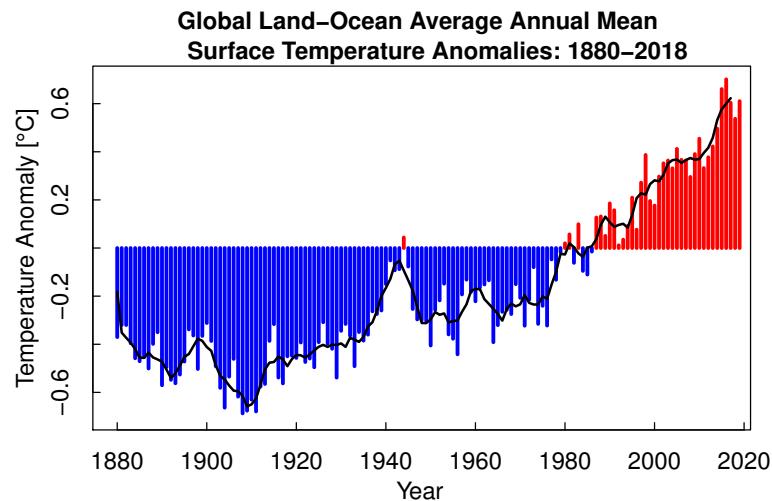
```
# Python plot Fig. 1.2: A staircase chart of data
# keyword arguments
kwargs = {'drawstyle' : 'steps'}
plt.plot(x, y, 'black', linewidth = 3, **kwargs);
plt.title("Global\u2014Land\u2014Ocean\u2014Average\u2014Annual\u2014Mean\u2014\n\u2014Surface\u2014Temperature\u2014Anomaly:\u20141880\u2014\u20142018", pad = 15)
plt.xlabel("Year", size = 25, labelpad = 20)
plt.ylabel("Temperature\u2014Anomaly\u2014[$\degree$C]", size = 25, labelpad = 20)
plt.show()
```

6046

### 6047 B.1.2.3 Plot a bar chart with colors

6048 Figure B.3 shows the anomaly size for each year by a color bar: red for positive  
 6049 anomalies and blue for negative anomalies. This bar chart style has an advantage of  
 6050 visualizing climate extremes, since these extremes may leave a distinct impression  
 6051 on viewers. The black curve is a 5-point moving average of the annual anomaly data,  
 6052 computed for each year by the mean of that year, together with two years before  
 6053 and two years after. The moving average smooths the high frequency fluctuations  
 6054 to reveal the long-term trends of temperature variations.

```
6055 # R plot Fig. 1.3: A color bar chart of data
6056 x <- NOAAtemp[,1]
6057 y <- NOAAtemp[,2]
6058 z <- rep(-99, length(x))
6059 # compute 5-point moving average
```



**Fig. B.3** Color bar chart of the 1880–2018 global average annual mean temperature anomalies.

```

6060 for (i in 3:length(x)-2) z[i] =
6061   mean(c(y[i-2],y[i-1],y[i],y[i+1],y[i+2]))
6062 n1 <- which(y>=0); x1 <- x[n1]; y1 <- y[n1]
6063 n2 <- which(y<0); x2 <- x[n2]; y2 <- y[n2]
6064 x3 <- x[2:length(x)-2]
6065 y3 <- z[2:length(x)-2]
6066 plot(x1, y1, type="h", #bars for data
6067   xlim = c(1880,2016), lwd=3,
6068   tck = 0.02, #tck>0 makes ticks inside the plot
6069   ylim = c(-0.7,0.7), xlab="Year", col="red",
6070   ylab = expression(paste(
6071     "Temperature", degree, "C")))
6072 main ="Global_Land-Ocean_Average_Annual_Mean
6073 Surface_Temperature_Anomalies_1880-2018",
6074 cex.lab = 1.2, cex.axis = 1.2)
6075 lines(x2, y2, type="h",
6076   lwd = 3, tck = -0.02, col = "blue")
6077 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\u2014Land\u2014Ocean\u2014Average\u2014Annual\u2014Mean\u2014\n\u2014Surface\u2014Temperature\u2014Anomaly:\u20141880\u2014\u20142018", 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\u2014Anomaly\u2014[$\degree$C]", size = 25, labelpad = 20)
plt.show()

```

6078

6079  
6080

### B.1.3 Statistical indices

6081 The commonly used basic statistical indices include mean, variance , standard de-  
 6082 viation, skewness, kurtosis, and quantiles. We first use R to calculate these indices  
 6083 for the global average annual mean temperature anomalies. Then we describe their  
 6084 mathematical formulas and interpret the numerical results.

```

6085 #R code for computing statistical indices
6086 setwd("/Users/ssten/climstats")
6087 NOAAtemp = read.table(
6088   "data/aravg.ann.land_ocean.90S.90N.v4.0.1.201907.txt",
6089   header=FALSE)
6090 temp2018=NOAAtemp[1:139,2] #use the temp data up to 2018
6091 head(temp2018) #show the first six values
6092 #[1] -0.370221 -0.319993 -0.320088 -0.396044 -0.458355 -0.470374
6093 mean(temp2018) #mean
6094 #[1] -0.1858632
6095 sd(temp2018) #standard deviation
6096 #[1] 0.324757
6097 var(temp2018) #variance
6098 #[1] 0.1054671

```

```

6099 library(e1071)
6100 #This R library is needed to compute the following parameters
6101 #install.packages("e1071") #if it is not in your computer
6102 skewness(temp2018)
6103 #[1] 0.7742704
6104 kurtosis(temp2018)
6105 #[1] -0.2619131
6106 median(temp2018)
6107 #[1] -0.274434
6108 quantile(temp2018,probs= c(0.05, 0.25, 0.75, 0.95))
6109 #   5%      25%     75%    95%
6110 # -0.5764861 -0.4119770  0.0155245  0.4132383

```

6111 We use  $x = \{x_1, x_2, \dots, x_n\}$  to denote the sampling data for a time series. The  
6112 statistical indices computed above by R are based on the following mathematical  
6113 formulas for mean, variance, standard deviation, skewness, and kurtosis:

$$\text{Mean: } \mu(x) = \frac{1}{n} \sum_{k=1}^n x_k, \quad (\text{B.3})$$

$$\text{Variance by unbiased estimate: } \sigma^2(x) = \frac{1}{n-1} \sum_{k=1}^n (x_k - \mu(x))^2, \quad (\text{B.4})$$

$$\text{Standard deviation: } \sigma(x) = (\sigma^2(x))^{1/2}, \quad (\text{B.5})$$

$$\text{Skewness: } \gamma_3(x) = \frac{1}{n} \sum_{k=1}^n \left( \frac{x_k - \mu(x)}{\sigma} \right)^3, \quad (\text{B.6})$$

$$\text{Kurtosis: } \gamma_4(x) = \frac{1}{n} \sum_{k=1}^n \left( \frac{x_k - \mu(x)}{\sigma} \right)^4 - 3. \quad (\text{B.7})$$

6114 A *quantile* cuts a sorted sequence of data. For example, the 25th quantile, also  
6115 called 25th percentile or 25% quantile, is the value that has 25% of the sorted data  
6116 smaller than this value and 75% larger than this value. The 50th percentile is also  
6117 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()

```

6118

6119 The meaning of these indices may be explained as follows. The mean is the simple  
 6120 average of samples. The variance of climate data reflects the strength of variations  
 6121 of a climate system and has units of the square of the data entries, such as  $[\text{ }^{\circ}\text{C}]^2$ .  
 6122 You may have noticed the denominator  $n - 1$  instead of  $n$ , which is for an estimate  
 6123 of unbiased sample variance. The standard deviation describes the spread of the  
 6124 sample entries and has the same units as the data. A large standard deviation  
 6125 means that the samples have a broad spread.

6126 Skewness is a dimensionless quantity. It measures the asymmetry of sample data.  
 6127 Zero skewness means a symmetric distribution about the sample mean. For example,  
 6128 the skewness of a normal distribution is zero. Negative skewness denotes a skew to  
 6129 the left, meaning the existence of a long tail on the left side of the distribution.  
 6130 Positive skewness implies a long tail on the right side.

6131 The words “Kurtosis” and “kurtic” are Greek in origin and indicate peakedness.  
 6132 Kurtosis is also dimensionless and indicates the degree of peakedness of a probability  
 6133 distribution. The kurtosis of a normal distribution<sup>1</sup> is zero when 3 is subtracted as  
 6134 in Eq. (B.7). Positive kurtosis means a high peak at the mean, thus the distribution  
 6135 shape is slim and tall. This is referred to as leptokurtic. “Lepto” is Greek in origin  
 6136 and means thin or fine. Negative kurtosis indicates a low peak at the mean, thus  
 6137 the distribution shape is fat and short, referred to as platykurtic. “Platy” is also  
 6138 Greek in origin and means flat or broad.

6139 For the 139 years of the NOAAGlobalTemp global average annual mean tem-  
 6140 perature anomalies, mean is  $-0.1959^{\circ}\text{C}$ , which means that the 1880-2018 average is  
 6141 lower than the average during the climatology period 1971-2000. During the clima-

<sup>1</sup> Chapter 2 will have a detailed description of the normal distribution and other probabilistic distributions.

6142 tology period, the temperature anomaly average is approximately zero, and can be  
 6143 computed by the R command `mean(temp2018[92:121])`.

6144 The variance of the data in the 139 years from 1880 to 2018 is  $0.1055 [{}^{\circ}\text{C}]^2$ ,  
 6145 and the corresponding standard derivation is  $0.3248 [{}^{\circ}\text{C}]$ . The skewness is 0.7743,  
 6146 meaning skew to the right with a long tail on the right, thus with more extreme high  
 6147 temperatures than extreme low temperatures, as shown by Fig. B.4. The kurtosis  
 6148 is -0.2619, meaning the distribution is flatter than a normal distribution, as shown  
 6149 in the histogram Fig. B.4.

6150 The median is  $-0.2744 {}^{\circ}\text{C}$  and is a number characterizing a set of samples such  
 6151 that 50% of the sample values are less than the median, and another 50% are  
 6152 greater than the median. To find the median, sort the samples from the smallest to  
 6153 the largest. The median is then the sample value in the middle. If the number of  
 6154 the samples is even, then the median is equal to the mean of the two middle sample  
 6155 values.

6156 Quantiles are defined in a way similar to median by sorting. For example, 25-  
 6157 percentile (also called the 25th percentile)  $-0.4120 {}^{\circ}\text{C}$  is a value such that 25% of  
 6158 sample data are less than this value. By definition, 60-percentile is thus larger than  
 6159 50-percentile. Here, percentile is a description of quantile relative to 100. Obviously,  
 6160 100-percentile is the largest datum, and 0-percentile is the smallest one. Often, a  
 6161 box plot is used to show the typical quantiles (see Fig. B.5).

6162 The 50-percentile (or 50th percentile)  $-0.2744 {}^{\circ}\text{C}$  is equal to the median. If the  
 6163 distribution is symmetric, then the median is equal to mean. Otherwise these two  
 6164 quantities are not equal. If the skew is to the right, then the mean is on the right of  
 6165 the median: the mean is greater than the median. If the skew is to the left, then the  
 6166 mean is on the left of the median: the mean is less than the median. Our 139 years  
 6167 of temperature data are right skewed and have mean equal to  $-0.1959 {}^{\circ}\text{C}$ , greater  
 6168 than their median equal to  $-0.2969 {}^{\circ}\text{C}$ .

## 6169 B.2 Commonly used climate statistical plots

---

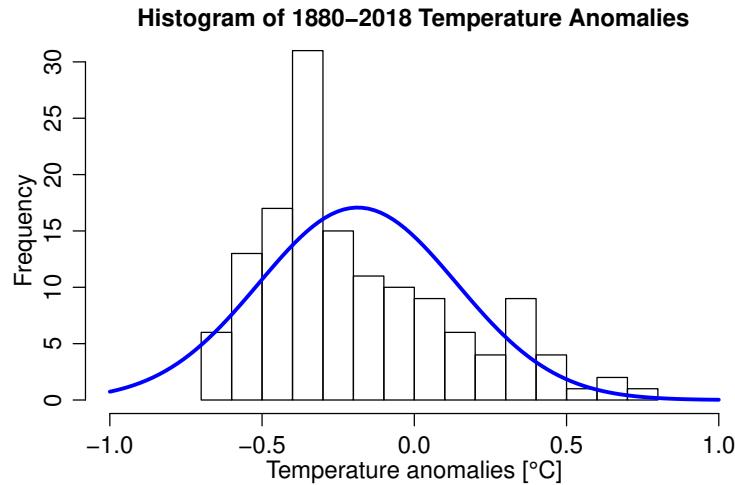
6170 We will use the 139 years of NOAAGlobalTemp temperature data and R to illustrate  
 6171 some commonly used statistical figures: histogram, boxplot, q-q plot, and linear  
 6172 regression trend line.

### 6174 B.2.1 Histogram of a set of data

---

6175 A histogram of the NOAAGlobalTemp global average annual mean temperature  
 6176 anomalies data from 1880-2018 is shown in Fig. B.4, which can be generated by the  
 6177 following R code.

```
6178 #R plot Fig. 1.4: Histogram and its fit
6179 par(mar=c(3.5,3.5,2.5,1), mgp=c(2,0.8,0))
6180 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.

```

6182 main="Histogram of 1880-2018 Temperature Anomalies",
6183 xlab=expression(paste(
6184   "Temperature anomalies [", degree, "C]")),
6185 xlim=c(-1,1), ylim=c(0,30),
6186 breaks=10, cex.lab=1.2, cex.axis=1.2)
6187 xfit <- seq(-1, 1, length=100)
6188 areat <- sum((h$counts)*diff(h$breaks[1:2]))#Normalization area
6189 yfit <- areat*dnorm(xfit,
6190   mean=mean(NOAAtemp[1:139,2]),
6191   sd=sd(NOAAtemp[1:139,2]))
6192 #Plot the normal fit on the histogram
6193 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 of 1880-2018 Temperature Anomalies",
          pad = 20)
plt.xlabel("Temperature Anomalies [$\degree$C]",
           size = 25, labelpad = 20)
plt.ylabel("Frequency", size = 25, labelpad = 20)
plt.show()

```

6194

6195     The shape of the histogram agrees with the characteristics predicted by the sta-  
 6196     tistical indices in the previous sub-section:

- 6197     (i) The distribution is asymmetric and skewed to the right with skewness equal to  
 6198        0.7743.
- 6199     (ii) The distribution is platykurtic with a kurtosis equal to -0.2619, i.e., it is flatter  
 6200        than the standard normal distribution indicated by the blue curve.

6201  
6202

## B.2.2 Box plot

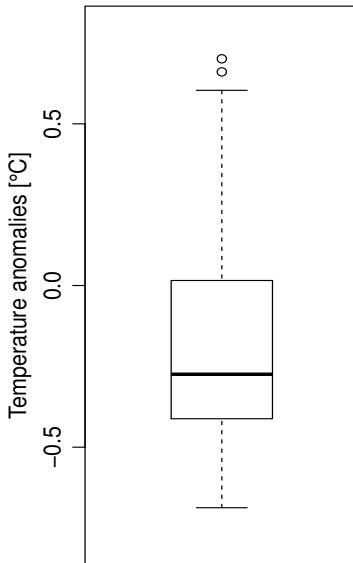
6203     Figure B.5 is the box plot of the 1880-2018 NOAAGlobalTemp global average annual  
 6204     mean temperature anomaly data, and can be made from the following R command

```

#R plot Fig. 1.5: Box plot
boxplot(NOAAtemp[1:139, 2], ylim = c(-0.8, 0.8),
        ylab=expression(paste(
            "Temperature anomalies [", degree, "C]")),
        width=NULL, cex.lab=1.2, cex.axis=1.2)

```

6210     The rectangular box's mid line indicates the median, which is  $-0.2744^{\circ}\text{C}$ . The  
 6211     rectangular box's lower boundary is the first quartile, i.e., 25th percentile  $-0.4120^{\circ}\text{C}$ .  
 6212     The box's upper boundary is the third quartile, i.e., the 75th percentile  $0.0155^{\circ}\text{C}$ .  
 6213     The box's height is the third quartile minus the first quartile, and is called the  
 6214     interquartile range (IQR). The upper "whisker" is the third quartile plus 1.5 IQR.  
 6215     The lower whisker is supposed to be at the first quartile minus 1.5 IQR. However,  
 6216     this whisker would then be lower than the lower extreme. In this case, the lower  
 6217     whisker takes the value of the lower extreme, which is  $-0.6872^{\circ}\text{C}$ . The points outside  
 6218     of the two whiskers are considered outliers. Our dataset has two outliers, which are  
 6219      $0.6607$  and  $0.7011^{\circ}\text{C}$ , and are denoted by two small circles in the box plot. The  
 6220     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()
```

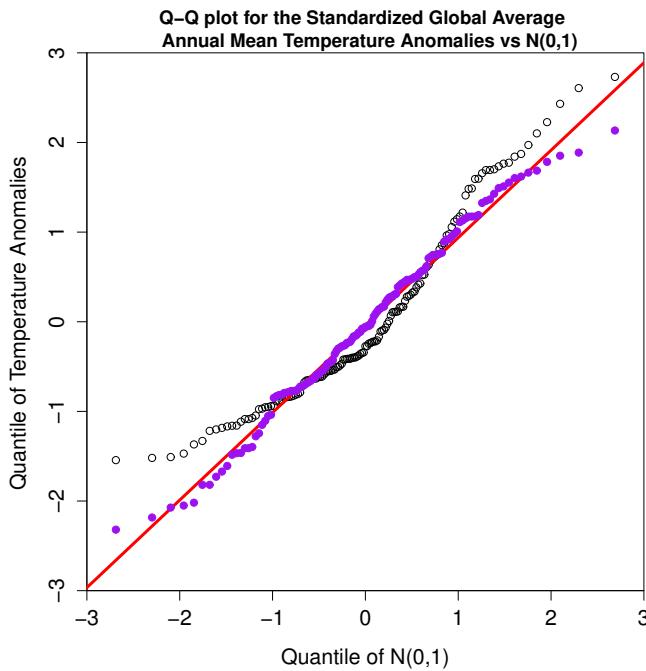
6221

6222  
6223

### B.2.3 Q-Q plot

6224 Figure B.6 shows Quantile-Quantile (Q-Q) plots, also denoted by q-q plots,  
6225 qq-plots, or QQ-plots.

6226 The function of a Q-Q plot is to compare the distribution of a given set of  
6227 data with a specific reference distribution, such as a standard normal distribution  
6228 with zero mean and standard deviation equal to one, denoted by  $N(0, 1)$ . A Q-Q  
6229 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-

6248   tion. However, we may also conclude that the distribution of these temperatures  
 6249   is not very far away from the normal distribution either, because the points on  
 6250   the Q-Q plot are not very far away from the distribution reference line, and also  
 6251   because even the simulated  $N(0, 1)$  points are noticeably off the reference line for  
 6252   the extremes.

6253   Figure B.6 can be generated by the following R code.

```

6254 # R plot Fig. 1.6: Q-Q plot for the standardized
6255 # global average annual mean temperature anomalies
6256 temp2018 <- NOAAtemp[1:139,2]
6257 tstand <- (temp2018 - mean(temp2018))/sd(temp2018)
6258 set.seed(101)
6259 qn <- rnorm(139) #simulate 139 points by N(0,1)
6260 qns <- sort(qn) # sort the points
6261 qq2 <- qqnorm(qns,col="blue", lwd = 2)
6262
6263 setEPS() #Automatically saves the eps file
6264 postscript("fig0106.eps", height=7, width=7)
6265 par(mar = c(4.5,5,2.5,1), xaxs = "i", yaxs = "i")
6266 qt = qnorm(tstand,
6267   main = "Q-Qplot for the Standardized Global Average
6268   Annual Mean Temperature Anomalies vs N(0,1)",
6269   ylab="Quantile of Temperature Anomalies",
6270   xlab="Quantile of N(0,1)",
6271   xlim=c(-3,3), ylim = c(-3,3),
6272     cex.lab = 1.3, cex.axis = 1.3)
6273 qqline(tstand, col = "red", lwd=3)
6274 points(qq2$x, qq2$y, pch = 19,
6275   col ="purple")
6276 dev.off()

```

6277   In the R code, we first standardize, also called normalize, the global average  
 6278   annual mean temperature data by subtracting the data mean and dividing by the  
 6279   data's standard deviation. Then, we use these 139 years of standardized global  
 6280   average annual mean temperature anomalies to generate a Q-Q plot, which is shown  
 6281   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\u2022plot\u2022for\u2022the\u2022Standardized\u2022Global\u2022\n\u2022Temperature\u2022Anomalies\u2022vs\u2022N(0,1)", pad = 20)
plt.xlabel("Quantile\u2022of\u2022N(0,1)", size = 25, labelpad = 20)
plt.ylabel("Quantile\u2022of\u2022Temperature\u2022Anomalies", size = 25)
plt.show()

```

6282

6283  
6284

## B.2.4 Plot a linear trend line

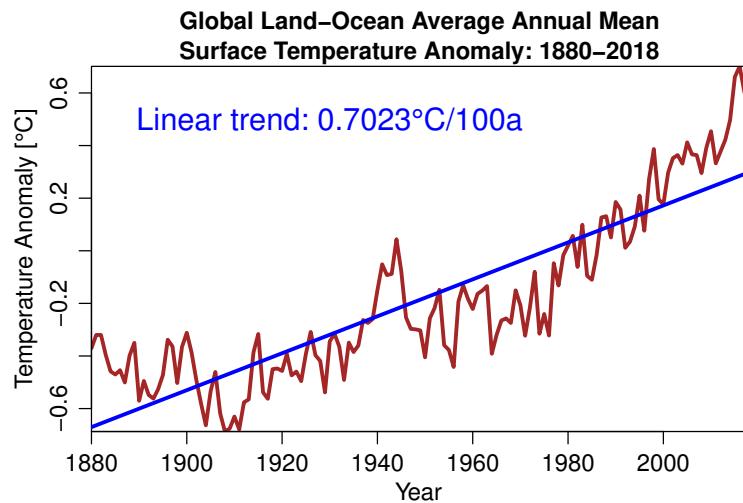
6285 Climate data analysis often involves plotting a linear trend line for time series data,  
 6286 such as the linear trend for the global average annual mean surface temperature  
 6287 anomalies, shown in Fig. B.7. The R code for plotting a linear trend line of data  
 6288 sequence  $y$  and time sequence  $t$  is `abline(y ~ t)`.

6289 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

```



**Fig. B.7** Linear trend line with the 1880–2018 global average annual mean surface temperature based on the NOAAGlobalTemp V4.0 dataset.

```

6307   text(1930, 0.5,
6308     expression(paste("Linear_trend:_0.7023",
6309                   "degree,C/100a")),
6310       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)

```

6311

### 6312 B.3 Read netCDF data file and plot spatial data maps

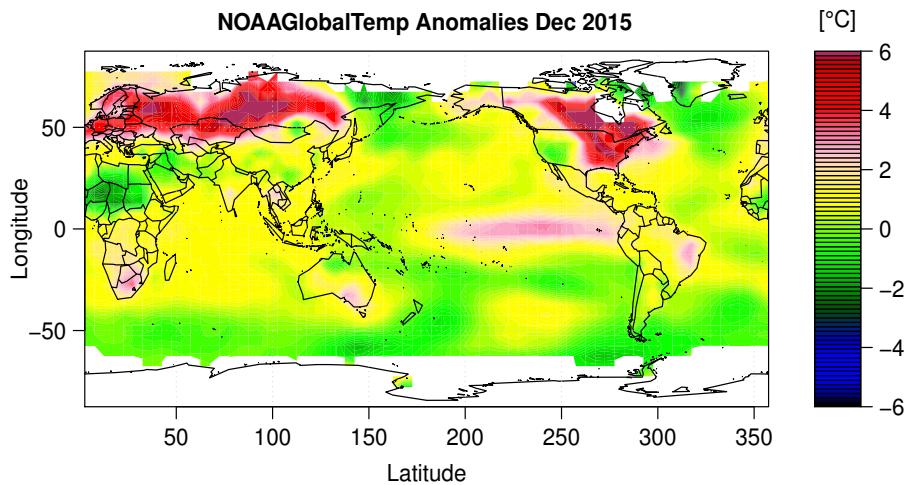
6313

#### 6314 B.3.1 Read netCDF data

6315 Climate data are at spatiotemporal points, such as at the grid points on the Earth's surface and at a sequence of time. NetCDF (Network Common Data Form) is a popular file format for modern climate data with spatial locations and temporal records. The gridded NOAAGlobalTemp data has a netCDF version, and can be downloaded from

6321 <https://www.esrl.noaa.gov/psd/data/gridded/data.noaaglobaltemp.html>

6322 The data are written into a 3D array, with 2D latitude-longitude for space, and 1D  
 6323 for time. R and Python can read and plot the netCDF data. We use the NOAA-  
 6324 GlobalTemp as an example to illustrate the netCDF data reading and plotting.  
 6325 Figure B.8 displays a temperature anomaly map for the entire globe for December  
 6326 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.

6327 Figure B.8 can be generated by the following computer code.

```
6328 # R read the netCDF data: NOAAGlobalTemp
6329 setwd("/Users/sshenn/climstats")
6330 #install.packages("ncdf4")
6331 library(ncdf4)
6332 nc = ncdf4::nc_open("data/air.mon.anom.nc")
6333 nc # describes details of the dataset
6334 Lat <- ncvar_get(nc, "lat")
6335 Lat # latitude data
6336 #[1] -87.5 -82.5 -77.5 -72.5 -67.5 -62.5
6337 Lon <- ncvar_get(nc, "lon")
```

```

6338 Lon # longitude data
6339 #[1] 2.5 7.5 12.5 17.5 22.5 27.5
6340 Time <- ncvar_get(nc, "time")
6341 head(Time) # time data in Julian days
6342 #[1] 29219 29250 29279 29310 29340 29371
6343 library(chron) # convert Julian date to calendar date
6344 nc$dim$time$units # .nc base time for conversion
6345 #[1] "days since 1800-1-1 00:00:0.0"
6346 month.day.year(29219,c(month = 1, day = 1, year = 1800))
6347 #1880-01-01 # the beginning time of the dataset
6348 tail(Time)
6349 #[1] 79988 80019 80047 80078 80108 80139
6350 month.day.year(80139,c(month = 1, day = 1, year = 1800))
6351 #2019-06-01 # the end time of the dataset
6352
6353 # extract anomaly data in (lon, lat, time) coordinates
6354 NOAAgridT <- ncvar_get(nc, "air")
6355 dim(NOAAgridT) # dimensions of the data array
6356 #[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/sshenn/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']
# convert 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)

```

6357

### B.3.2 Plot a spatial map of temperature

6358  
6359 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.

6360  
6361  
6362  
6363  
6364  
6365  
6366 The figure can be generated by the following R code.

```
6367 #R plot Fig. 1.8: Dec 2015 surface temp anomalies map
6368 library(maps)# requires maps package
6369 mapmat=NOAAgridT[,1632]
6370 # Julian date time 1632 corresponds to Dec 2015
6371 mapmat=pmax(pmin(mapmat,6),-6) # put values in [-6, 6]
6372 int=seq(-6, 6, length.out=81)
6373 rgb.palette=colorRampPalette(c('black','blue',
6374   'darkgreen', 'green', 'yellow','pink','red','maroon'),
6375   interpolate='spline')
6376 par(mar=c(3.5, 4, 2.5, 1), mgp=c(2.3, 0.8, 0))
6377 filled.contour(Lon, Lat, mapmat,
6378   color.palette=rgb.palette, levels=int,
6379   plot.title=title(main="NOAAGlobalTemp\u2014Anomalies\u2014Dec\u20142015",
6380     xlab="Latitude", ylab="Longitude", cex.lab=1.2),
6381   plot.axes={axis(1, cex.axis=1.2, las=1);
6382     axis(2, cex.axis=1.2, las=2); map('world2', add=TRUE); grid()},
6383   key.title=title(main=expression(paste("[", degree, "C]"))),
6384   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('[$\degree$C]', size= 17, pad = 10)
cbar.ax.tick_params(labelsize = 15)
# add plot title
plt.title('NOAA Global Temp Anomalies Dec 2015',
           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()

```

6385

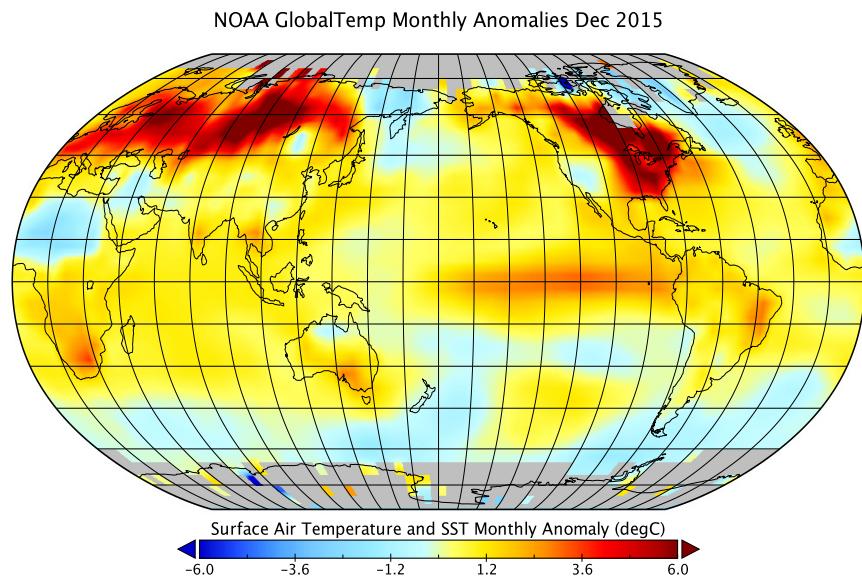
6386

6387

### B.3.3 Panoply plot of a spatial map of temperature

6388 You can also use the Panoply software package to plot the map (See Fig. B.9).  
 6389 This is a very powerful data visualization tool developed by NASA specifically for  
 6390 displaying netCDF files. The software package is free and can be downloaded from  
 6391 [www.giss.nasa.gov/tools/panoply](http://www.giss.nasa.gov/tools/panoply).

6392 To make a Panoply plot, open Panoply and choose Open from the File menu.  
 6393 Open dropdown which will allow you to go to the right directory to find the netCDF  
 6394 file you wish to plot. In our case, the file is `air.mon.anom.nc`. Choose the climate  
 6395 parameter `air`. Click on Create Plot. A map will show up. Then you have many  
 6396 choices to modify the map, ranging from the data `Array`, `Scale`, and `Map`, ...,  
 6397 and finally produce the figure. You can then tune the figure by choosing different  
 6398 graphics parameters underneath the figure, such as `Array(s)` to choose which



**Fig. B.9** A Panoply plot of Robinson projection map for the surface temperature anomalies of December 2015.

month to plot, **Map** to choose the map projection types and the map layout options, and **Labels** to type proper captions and labels.

To learn more about Panoply, please use an online Panoply tutorial, such as the NASA Panoply help page <https://www.giss.nasa.gov/tools/panoply/help/>.

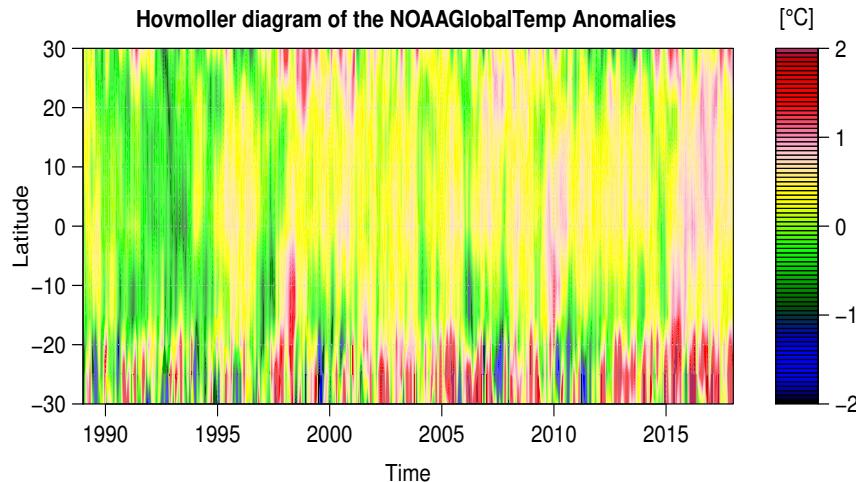
Compared with the R or Python map Fig. B.5, the visualization effect of the Panoply map seems more appealing. However, R and Python have an advantage of flexibility and can deal with all kinds of data. For example, the Plotly graphing library in R <https://plotly.com/r/> and in Python <https://plotly.com/python/> can even make interactive and 3D graphics. You may find some high quality figures from the Intergovernmental Panel on Climate Change (IPCC) report (2021) [www.ipcc.ch](http://www.ipcc.ch) and reproduce them using the computing tools described here.

## B.4 1D-space-1D-time data and Hovmöller diagram

A very useful climate data visualization technique is the Hovmöller diagram. It displays how a climate variable varies with respect to time along a given line section of latitude, longitude or altitude. It usually has the abscissa for time and ordinate for the line section. A Hovmöller diagram can conveniently show time evolution of a spatial pattern, such as wave motion from south to north or from west to east.

Figure B.10 is a Hovmöller diagram for the sea surface temperature (SST) anomalies at a longitude equal to  $240^{\circ}$ , i.e.,  $120^{\circ}\text{W}$ , in a latitude interval  $[30^{\circ}\text{S}, 30^{\circ}\text{N}]$ ,

6419 with time range from January 1989 to December 2018. When the red strips become strong from the south to north, a strong El Niño occurs, such as those in the  
 6420 1997-1998 and 2015-2016 winters.  
 6421



**Fig. B.10** Hovmöller diagram for the gridded NOAAGlobalTemp monthly anomalies at longitude  $120^{\circ}\text{W}$  and a latitude interval  $[30^{\circ}\text{S}, 30^{\circ}\text{N}]$ .

6422 The Hovmöller diagram Fig. B.10 may be plotted by the following R code.

```
6423 #R plot Fig. 1.10: Hovmoller diagram
6424 library(maps)
6425 mapmat=NOAAgridT[30,12:24,1309:1668]
6426 #Longitude= 240 deg, Lat =[-30 30] deg
6427 #Time=Jan 1989-Dec 2018: 30 years
6428 mapmat=pmax(pmin(mapmat,2),-2) # put values in [-2,2]
6429 par(mar=c(4,5,3,0))
6430 int=seq(-2,2,length.out=81)
6431 rgb.palette=colorRampPalette(c('black','blue',
6432   'darkgreen','green','yellow','pink','red','maroon'),
6433   interpolate='spline')
6434 par(mar=c(3.5,3.5,2.5,1), mgp=c(2.4, 0.8, 0))
6435 x = seq(1989, 2018, len=360)
6436 y = seq(-30, 30, by=5)
6437 filled.contour(x, y, t(mapmat),
6438   color.palette=rgb.palette, levels=int,
6439   plot.title=title(main=
6440     "Hovmoller diagram of the NOAAGlobalTemp Anomalies",
6441     xlab="Time", ylab="Latitude", cex.lab=1.2),
6442   plot.axes={axis(1, cex.axis=1.2);
6443     axis(2, cex.axis=1.2);
6444     map('world2', add=TRUE); grid()},
6445   key.title=title(main =
6446     expression(paste("[", degree, "C]"))),
6447   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\u2022diagram\u2022of\u2022the\u2022NOAA\u2022GlobalTemp\u2022Anomalies",
          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])

```

6448

6449

## B.5 4D netCDF file and its map plotting

6450

6451 The 4D climate data means 3D spatial dimensions and 1D time. For example, the  
 6452 NCEP Global Ocean Data Assimilation System (GODAS) monthly water temper-  
 6453 ature data are at 40 depth levels ranging from 5 meters to 4478 meters and at 1/3  
 6454 degree latitude by 1 degree longitude horizontal resolution and are from January  
 6455 1980. The NOAA-CIRES 20th Century Reanalysis (20CR) monthly air tempera-  
 6456 ture data are at 24 different pressure levels ranging from 1000 mb to 10 mb and at  
 6457 2-degree latitude and longitude horizontal resolution and are from January 1851.

6458 GODAS data can be downloaded from NOAA ESRL

6459 [www.esrl.noaa.gov/psd/data/gridded/data.godas.html](http://www.esrl.noaa.gov/psd/data/gridded/data.godas.html)

6460 The data for each year is a netCDF file and is about 140MB. The following R code  
 6461 can read the GODAS 2015 data into R.

```

6462 # R read a 4D netCDF file: lon, lat, level, time
6463 setwd("/Users/sshen/climstats")
6464 library(ncdf4)
6465 # read GODAS data 1-by-1 deg, 40 levels, Jan-Dec 2015
6466 nc=ncdf4::nc_open("data/godas2015.nc")
6467 nc
6468 Lat <- ncvar_get(nc, "lat")
6469 Lon <- ncvar_get(nc, "lon")
6470 Level <- ncvar_get(nc, "level")
6471 Time <- ncvar_get(nc, "time")
6472 head(Time)
6473 #[1] 78527 78558 78586 78617 78647 78678
6474 library(chron)
6475 month.day.year(78527,c(month = 1, day = 1, year = 1800))
6476 # 2015-01-01
6477 # potential temperature pottmp[lon, lat, level, time]
6478 godasT <- ncvar_get(nc, "pottmp")
6479 dim(godasT)
6480 #[1] 360 418 40 12,
6481 #i.e., 360 lon, 418 lat, 40 levels, 12 months=2015
6482 t(godasT[246:250, 209:210, 2, 12])
6483 #Dec level 2 (15-meter depth) water temperature [K] of
6484 #a few grid boxes over the eastern tropical Pacific
6485 # [,1] [,2] [,3] [,4] [,5]
6486 #[1,] 300.0655 299.9831 299.8793 299.7771 299.6641
6487 #[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)

```

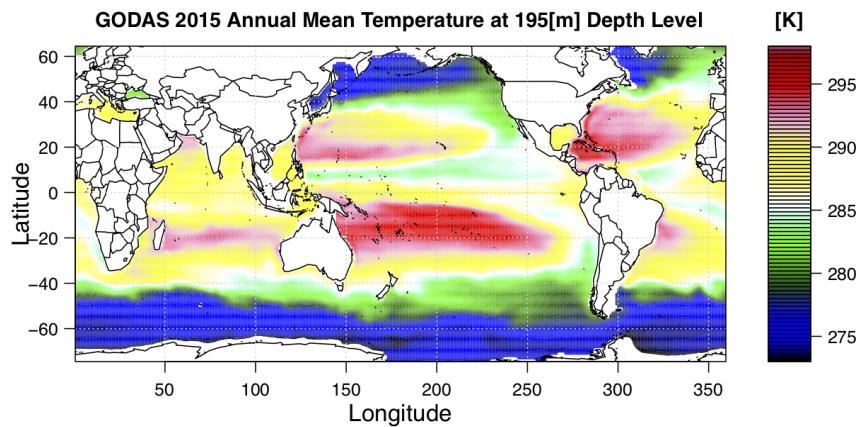
6488 Figure B.11 shows the 2015 annual mean water temperature at 195 meters depth  
 6489 based on the GODAS data. At this depth level, the equatorial upwelling appears:  
 6490 The deep ocean cooler water in the equatorial region upwells and makes the equa-  
 6491 torial water cool. The equatorial water is not the hottest anymore at this level, and  
 6492 is cooler than the water in some subtropical regions as shown in Fig. B.11.  
 6493

6494 Figure B.11 can be generated by the following R code.

```

6495 # R plot Fig. 1.11: The ocean potential temperature
6496 # the 20th layer from surface: 195 meters depth
6497 # compute 2015 annual mean temeprature at 20th layer

```



**Fig. B.11** The 2015 annual mean water temperature at 195 meters depth based on GODAS data.

```

6498 library(maps)
6499 climmat=matrix(0,nrow=360,ncol=418)
6500 sdmat=matrix(0,nrow=360,ncol=418)
6501 Jmon<-1:12
6502 for (i in 1:360){
6503   for (j in 1:418){
6504     climmat[i,j] = mean(godasT[i,j,20,Jmon]);
6505     sdmat[i,j]=sd(godasT[i,j,20,Jmon])
6506   }
6507 }
6508 int=seq(273,298,length.out=81)
6509 rgb.palette=colorRampPalette(c('black','blue',
6510   'darkgreen','green','white','yellow',
6511   'pink','red','maroon'), interpolate='spline')
6512 par(mar=c(3.5, 3.5, 2.5, 0), mgp=c(2, 0.8, 0))
6513 filled.contour(Lon, Lat, climmat,
6514   color.palette=rgb.palette, levels=int,
6515   plot.title=title(main=
6516   "GODAS 2015 Annual Mean Temperature at 195[m] Depth Level",
6517   xlab="Longitude", ylab="Latitude",
6518   cex.lab=1.3, cex.axis=1.3),
6519   plot.axes={axis(1); axis(2); map('world2', add=TRUE); grid()},
6520   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_2015_Annual_Mean_Temperature\n"
          "at_195[m]_Depth_Level",
          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])

```

6521

6522

6523

## B.6 Paraview and 4DVD

6524 Besides using R, Python and Panoply to plot climate data, other software packages  
 6525 may also be used to visualize data for some specific purposes, such as Paraview for  
 6526 3D visualization and 4DVD for fast climate diagnostics and data delivery.

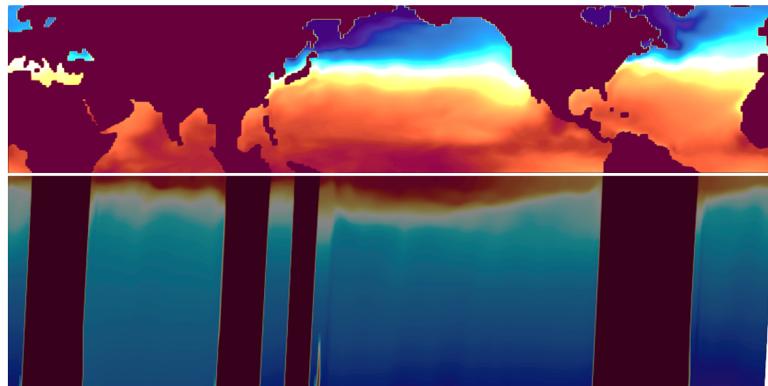
6527

6528

### B.6.1 Paraview

6529 Paraview is an open-source data visualization software package, available at  
 6530 [www.paraview.org/download](http://www.paraview.org/download). There are many online tutorials, such as  
 6531 [www.paraview.org/tutorials](http://www.paraview.org/tutorials).

6532 You may use the software ParaView to plot the GODAS data in a 3D view as  
 6533 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.

### B.6.2 4DVD

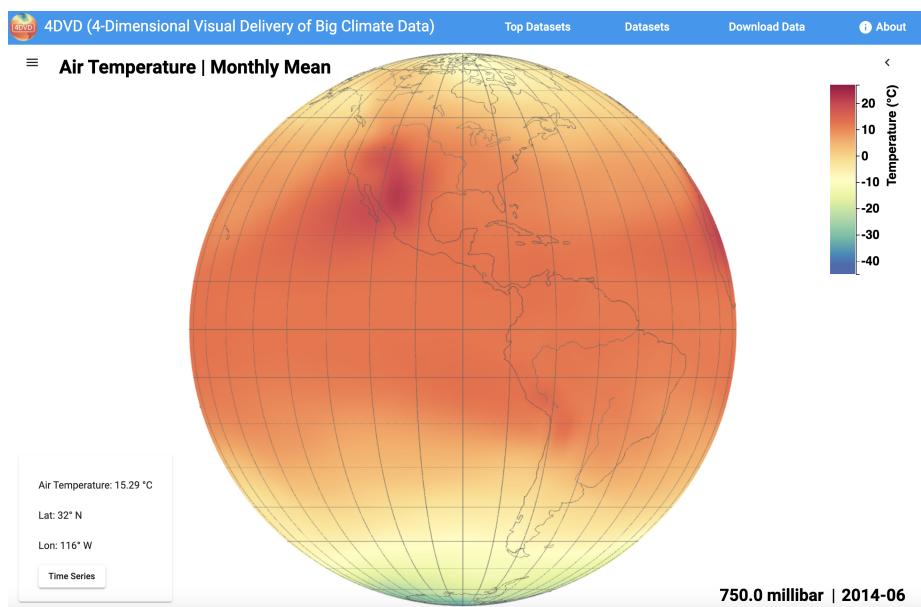
6534 4DVD (4-dimensional visual delivery of big climate data) is a fast data visualization  
 6535 and delivery software system [www.4dvd.org](http://www.4dvd.org). It optimally harnesses the power of  
 6536 distributed computing, database and data storage to allow a large number of general  
 6537 public users to quickly visualize climate data. For example, teachers and students  
 6538 can use 4DVD to visualize climate model data in classrooms and download the  
 6539 visualized data instantly. The 4DVD website has a tutorial for users.

6540 Here, we provide an example of 4DVD visualization of the NOAA-CIRES 20th  
 6541 Century Reanalysis climate model data for atmosphere. 4DVD can display a tem-  
 6542 perature map at a given time and given pressure level as shown in Fig. B.13 for  
 6543 January 1851 and 750 millibar pressure level, or one can obtain a time series of  
 6544 the monthly air temperature for a specific grid point from January 1851. In fact,  
 6545 4DVD can show multiple time series for the same latitude-longitude location but  
 6546 at different pressure levels. The 4DVD not only allows a user to visualize the data,  
 6547 but also download the data for the figure shown in the 4DVD system. In this sense,  
 6548 4DVD is like a machine that plays data, while the regular DVD player machine,  
 6549 popular for about 30 years since 1980s, play DVD discs for music and movies.

### B.6.3 Other online climate data visualization tools

6552 Besides R, Python, Panoply, ParaView, and 4DVD, there are many other data  
 6553 visualization and delivery software systems. A few popular and free ones are listed  
 6554 as follows.

6555 Nullschool <https://nullschool.net> is a beautiful data visualizer of wind, ocean  
 6556 flows, and many climate parameters. It is supported by the data from a global



**Fig. B.13** The 4DVD screenshot for the NOAA-CIRES 20th Century Reanalysis temperature at 750 millibar height level for June 2014.

numerical weather prediction system, named the Global Forecast System (GFS) run by the United States' National Weather Service (NWS).

Ventusky [www.ventusky.com](http://www.ventusky.com) has both website and smartphone app. It has an attractive and user-friendly interface that allows users to get digital weather information instantly around the globe.

Climate Reanalyzer [climatereanalyzer.org](http://climatereanalyzer.org) is a comprehensive tool for climate data plotting and download. It has a user-friendly interface for reanalysis and historical station data. The data can be exported in either CSV (comma-separated values) format or JSON (JavaScript object notation) format.

Google Earth Engine [earthengine.google.com](http://earthengine.google.com) provides visualization tools together with a huge multi-petabytes storage of climate data. Its modeling and satellite data sources are from multiple countries.

Giovanni [giovanni.gsfc.nasa.gov](http://giovanni.gsfc.nasa.gov) is an online climate data plotting tool with an interface. It allows to download the plotted figures in different formats, such as png. It is supported by various kinds of NASA climate datasets.

Climate Engine [climateengine.org](http://climateengine.org) is a web application for plotting climate and remote sensing data. Similar to Giovanni, it also has a tabular interface for a user to customized her data and maps.

NOAA Climate at a Glance [www.ncdc.noaa.gov/cag](http://www.ncdc.noaa.gov/cag) is a data visualization tool mainly for visualizing the observed climate data over the United States. It has functions for both spatial maps and historical time series.

Web-based Reanalyses Intercomparison Tools (WRIT) [psl.noaa.gov/data/writ](http://psl.noaa.gov/data/writ).

6581 Similar to Giovanni, WRIT also has an interface table for a user to enter plot parameters.  
 6582 The plot (in postscript format) and its data (in netCDF format) can be  
 6583 downloaded. WRIT is designed for the data from climate reanalysis models.

#### 6584 **B.6.4 Use ChatGPT as a study assistant**

6586 ChatGPT (Generative Pre-trained Transformer) is an artificial intelligence (AI)  
 6587 system that can help us develop computer codes, write essays, draft course syllabi,  
 6588 and more. For example, in your box of the ChatGPT interface, you can enter “Write  
 6589 a Python code to plot the cosine function as a blue thick dashed curve.” ChatGPT  
 6590 will give you the following Python code:

```
import numpy as np
import matplotlib.pyplot as plt

x = np.linspace(-np.pi, np.pi, 100)
y = np.cos(x)

plt.plot(x, y, linestyle='--', linewidth=3, color='blue')
plt.xlabel('x')
plt.ylabel('cos(x)')
plt.title('Plot of the cosine function')
plt.show()
```

6592 You can copy the above code and paste it to your Jupyter Notebook cell, then run  
 6593 the code to produce the curve of the cosine function in the interval  $[-\pi, \pi]$ . With  
 6594 this sample and according to your needs, you can modify the code for different  
 6595 functions, intervals, colors, thickness, title, axis labels, and more.

6596 You can also ask ChatGPT to work on a dataset, such as “Write an R code  
 6597 to plot the NOAAGlobalTempts.csv data and its trend line.” ChatGPT will give  
 6598 you an R code. You copy and paste the code to your RStudio. Set your R code in  
 6599 the correct directory that has the NOAAGlobalTemp annual time series data file,  
 6600 NOAAGlobalTempts.csv. Run your R code to plot the data and the linear trend  
 6601 line. Of course, you can modify the code according to your requirements. The R  
 6602 code is as follows.

```
6603 setwd('/Users/sshenn/climmath/data')
6604 library(tidyverse)
6605
6606 # Load the data
6607 data <- read.csv("NOAAGlobalTempts.csv")
6608
6609 # Plot the data
6610 ggplot(data, aes(x = Year, y = Tempts)) +
  geom_line() +
  geom_smooth(method = "lm", se = FALSE) +
  ggtitle("NOAA Global Temperature Data") +
  xlab("Year") +
  ylab("Temperature (Celsius)")
```

6616 Another example is that you enter “Write an essay on global warming.” ChatGPT  
6617 will write an essay for you. Most likely, it does not reflect your idea, but it will  
6618 give you a few tips, such as the definition of global warming, evidences and data,  
6619 consequences, main causes, how to slow down the warming, and conclusion. Again,  
6620 you can use this ChatGPT essay as a sample and modify the essay using your own  
6621 data and ideas, such as green technologies, solar farms, microclimate resources for  
6622 smart cities, trend of the oil business, lifestyles, AI in self-driving electric cars, and  
6623 education.

6624 ChatGPT was released by the OpenAI Lab in November 2022. You may use this  
6625 tool as your assistant of learning and working. It is only your assistant. It is not  
6626 you! You can efficiently use it to give you hints and inspire your ideas. You still  
6627 have to produce your own work.

## 6628 **B.7 Chapter summary**

---

6630 This chapter has provided a brief introduction to useful statistical concepts and  
6631 methods for climate science and has included the following materials.

6632 (i) Formulas to compute the most commonly used statistical indices:

- 6633 • Mean as a simple average, median as a datum whose value is in the  
6634 middle of the sorted data sequence, i.e., the median is larger than 50%  
6635 of the data and smaller than the remaining 50%,
- 6636 • Standard deviation as a measure of the width of the probability distri-  
6637 bution,
- 6638 • Variance as the square of the standard deviation,
- 6639 • Skewness as a measure of the degree of asymmetry in distribution, and
- 6640 • Kurtosis as a measure of the peakedness of the data distribution com-  
6641 pared to that of the normal distribution.

6642 (ii) The commonly used statistical plots:

- 6643 • Histogram for displaying the probability distribution of data,
- 6644 • Linear regression line for providing a linear model for data,
- 6645 • Box plot for quantifying the probability distribution of data, and
- 6646 • Q-Q plot for checking whether the data are normally distributed.

6647 (iii) Read and plot the netCDF file for plotting a 2D map by R and Python.  
6648 Other data visualization tools, such as Panoply, Paraview, 4DVD, Plotly,  
6649 and Nullschool, were briefly introduced. The online tools like 4DVD and  
6650 Nullschool may be used for classroom teaching and learning.

6651 Computer codes and climate data examples are given to demonstrate these con-  
6652 cepts and the use of the relevant tools and formulas. With this background plus  
6653 some R or Python programming skill, you will have sufficient knowledge to meet  
6654 the needs of the basic statistical analysis for climate data.

6655

## References and Further Readings

6656 [1] Hennemuth, B., Bender, S., K. Bulow, N. Dreier, E. Keup-Thiel, O. Kruger,  
6657 C. Mudersbach, C. Radermacher, and R. Schoetter, 2013: *Statistical methods*  
6658 *for the analysis of simulated and observed climate data, applied in projects and*  
6659 *institutions dealing with climate change impact and adaptation.* CSC Report 13,  
6660 Climate Service Center, Germany, 135pp. URL:

6661 [http://www.climate-service-center.de/imperia/md/content/csc/  
6662 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.

6663

6664 [2] IPCC, 2021: AR6 Climate Change 2021: The Physical Science Basis. Contribution  
6665 of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change [Masson-Delmotte, V., P. Zhai, A. Pirani, S.  
6666 L. Connors, C. Pan, S. Berger, N. Caud, Y. Chen, L. Goldfarb, M. I. Gomis, M.  
6667 Huang, K. Leitzell, E. Lonnoy, J. B. R. Matthews, T. K. Maycock, T. Waterfield,  
6668 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

6670

6671 [3] NCEI, 2021: *Anomalies vs. Temperature* Last accessed on 12 May 2021.

6672 [https://www.ncdc.noaa.gov/monitoring-references/dyk/  
6673 anomalies-vs-temperature](https://www.ncdc.noaa.gov/monitoring-references/dyk/anomalies-vs-temperature)

6674

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.

6675

- 6676 [4] Shen, S.S.P, and R.C.J. Somerville, 2019: Climate Mathematics: Theory and  
 6677 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.climate-mathematics.org](http://www.climate-mathematics.org).

- 6678  
 6679 [5] Smith, T.M., R.W. Reynolds, T.C. Peterson, and J. Lawrimore, 2008: Improvements to NOAA's historical merged landocean surface temperatures analysis  
 6680 (18802006). Journal of Climate, 21, 22832296, doi:10.1175/2007JCLI2100.1.

These authors are among the main contributors who reconstructed the sea surface temperature (SST). They began their endeavor in the early 1990s.

- 6681  
 6682 [6] Vose, R.S., D. Arndt, V.F. Banzon, D.R. Easterling, B. Gleason, B. Huang,  
 6683 E. Kearns, J.H. Lawrimore, M.J. Menne, T.C. Peterson, R.W. Reynolds, T.M.  
 6684 Smith, C.N. Williams, Jr., and D.L. Wuertz, 2012: NOAA's merged land-ocean  
 6685 surface temperature analysis. Bulletin of the American Meteorological Society,  
 6686 93, 16771685.

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.

## Exercises

- 6689  
 6690 **B.1** The NOAA Merged Land Ocean Global Surface Temperature Analysis (NOAA-  
 6691 GlobalTemp) V5 includes the global land-spatial average annual mean tem-  
 6692 perature anomalies as shown below

6694	1880	-0.843351	0.031336	0.009789	0.000850	0.020698
6695	1881	-0.778600	0.031363	0.009789	0.000877	0.020698
6696	1882	-0.802413	0.031384	0.009789	0.000897	0.020698
6697	.....					

- 6698      The first column is for time in years, and the second is the temperature anomalies  
6699      in [Kelvin]. Columns 3-6 are data errors. This data file for the land tem-  
6700      perature can be downloaded from the NOAAGlobalTemp website  
6701      [www.ncei.noaa.gov/data/noaa-global-surface-temperature/v5/access/timeseries](http://www.ncei.noaa.gov/data/noaa-global-surface-temperature/v5/access/timeseries)  
6702      The data file named `aravg.ann.land.90S.90N.v5.0.0.202104.asc.txt` is  
6703      also included in `data.zip` that can be downloaded from the book website  
6704      [www.climatestatistics.org](http://www.climatestatistics.org)  
6705      (a) Plot the anomalies against time from 1880 to 2020 using the point-line  
6706      graph like Fig. B.1.  
6707      (b) Plot the anomalies against time from 1880 to 2020 using the staircase  
6708      chart like Fig. B.2.  
6709      (c) Plot the anomalies against time from 1880 to 2020 using the color bar  
6710      chart like Fig. B.3.
- 6711 **B.2** NOAAGlobalTemp V5 also has the global **ocean**-spatial average annual mean  
6712      temperature anomalies contained in a data file named  
6713      `aravg.ann.locean.90S.90N.v5.0.0.202104.asc.txt`  
6714      For this dataset, please plot the anomaly data in the same three styles (a),  
6715      (b) and (c) as in the previous problem. You may download the data file from  
6716      the NOAAGlobalTemp V5 website or from `data.zip` for this book.
- 6717 **B.3** NOAAGlobalTemp V5 also has the global land-spatial average monthly mean  
6718      temperature anomalies from January 1880 to April 2021. The data file is  
6719      named  
6720      `aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt`  
6721      For this dataset, please plot the January anomaly data from January 1880 to  
6722      January 2021 in the same three styles (a) - (c) as in the previous problem.
- 6723 **B.4** For the monthly global land data of NOAAGlobalTemp V5 in the above prob-  
6724      lem,  
6725      (a) For the January data from 1880 to 2020, compute the following statistical  
6726      indices: Mean, standard deviation, variance, skewness, kurtosis, max, 75%-  
6727      percentile, median, 25%-percentile, and min.  
6728      (b) Do the same for February, March, ..., December.  
6729      (c) Aggregate all the results in (a) and (b) into a  $10 \times 12$  matrix with the 10  
6730      statistical indices in rows and the 12 months in columns. Add proper column  
6731      names, such as `Jan Feb Mar ...`, and also row names.  
6732      (d) Use 100-300 words to describe the differences of the statistical indices for  
6733      different months.
- 6734 **B.5** For the monthly data CRUTEM4 (Climate Research Unit surface air tempera-  
6735      ture anomalies) in 4DVD, download the historical time series data for January  
6736      of a location of your interest, and compute the following statistical indices:  
6737      Mean, standard deviation, variance, skewness, kurtosis, max, 75%-percentile,  
6738      median, 25%-percentile, and min.
- 6739 **B.6** Do the same as the previous problem but for July. Comment on the differences  
6740      between the statistical indices in January and July.
- 6741 **B.7** (a) Plot a histogram of the January global land temperature anomalies using

- 6742       the data file in the above problem, i.e.,  
 6743       **aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt**  
 6744       (b) Do the same for July.  
 6745       (c) Use 100-200 words to describe the comparison of the two histograms.
- 6746 **B.8** (a) Plot a box plot for the January global land temperature anomalies  
 6747       **aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt**  
 6748       (b) Do the same but for the July data.  
 6749       (c) Put the two box plots next to each other in the same figure.  
 6750       (d) Use 100-200 words to describe the comparison of the two box plots.
- 6751 **B.9** (a) Use the monthly data in the above problem to generate twelve (12) Q-  
 6752       Q plots relative to the standard normal distribution for every month from  
 6753       January to December.  
 6754       (b) From the 12 Q-Q plots, discuss which month's temperature anomalies are  
 6755       the closest to the normal distribution.
- 6756 **B.10** (a) Using the monthly data in the above problem  
 6757       **aravg.mon.land.90S.90N.v5.0.0.202104.asc.txt**  
 6758       compute the linear trend of January temperature anomalies from 1880 to  
 6759       2020. Output your result in the unit: [°C/century].  
 6760       (b) Repeat (a) for each of the other eleven months.  
 6761       (c) Generate a list for the twelve linear trends.  
 6762       (d) Use 100-200 words to discuss the numerical values of the list.
- 6763 **B.11** Plot the December 1997 surface temperature anomalies using the NOAA-  
 6764       GlobalTemp V5 data on a  $5^\circ \times 5^\circ$  grid in the netCDF format:
- 6765       **NOAAGlobalTemp\_v5.0.0\_gridded\_s188001\_e202104\_c20210509T133251.nc**
- 6766       The data can be downloaded from NOAAGlobalTemp V5 or extracted from  
 6767       **data.zip** for this book at the website [www.climatestatistics.org](http://www.climatestatistics.org)
- 6768 **B.12** Use Panoply to make the same December 1997 surface temperature anomalies  
 6769       plot as the above problem, but with the *Robinson* map projection.
- 6770 **B.13** Use Panoply to make the same December 1997 surface temperature anomalies  
 6771       plot as the above problem, but with the *Mercator* map projection.
- 6772 **B.14** Use Panoply to make the same December 1997 surface temperature anomalies  
 6773       plot as the above problem, but with the *Orthographic* map projection.
- 6774 **B.15** Use Panoply to make the same December 1997 surface temperature anomalies  
 6775       plot as the above problem, but with the *Stereographic (Two-Hemisphere)* map  
 6776       projection.
- 6777 **B.16** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp  
 6778       monthly anomalies at longitude  $150^\circ\text{W}$  and a latitude interval  $[50^\circ\text{S}, 50^\circ\text{N}]$   
 6779       from January 1989 to December 2018 (i.e., 240 months).
- 6780 **B.17** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp  
 6781       monthly anomalies at longitude  $140^\circ\text{W}$  and a latitude interval  $[40^\circ\text{S}, 40^\circ\text{N}]$   
 6782       from January 1971 to December 2000 (i.e., 360 months).

- 6783   **B.18** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp  
6784       monthly anomalies on the equator with longitude from 160°E to 90°W] from  
6785       January 1971 to December 2000 (i.e., 360 months).
- 6786   **B.19** Plot a Hovmöller diagram like Fig. B.10 for the gridded NOAAGlobalTemp  
6787       monthly anomalies at latitude 5°S with longitude from 160°E to 90°W] from  
6788       January 1971 to December 2000 (i.e., 360 months).
- 6789   **B.20** Use R or Python to plot four December 2015 potential temperature maps  
6790       at the depth layers 5, 25, 105 and 195 meters based on the GODAS data.  
6791       You can use the netCDF data file `godas2015.nc` in `data.zip` for this book  
6792       or download the netCDF GODAS data from a NOAA website, such as  
6793       <https://www.esrl.noaa.gov/psd/data/gridded/data.godas.html>
- 6794   **B.21** Use Panoply to plot the same maps as the previous problem but with the  
6795       Robinson projection.
- 6796   **B.22** Use Plotly in R or Python to plot four December 2015 potential temperature  
6797       maps at the depth layers 5, 25, 105 and 195 meters based on the GODAS data.  
6798       Try to place the four maps together to make a 3D visualization.
- 6799   **B.23** Plot four June 2015 potential temperature maps at the depth layers 5, 25,  
6800       105 and 195 meters using the same netCDF GODAS data file as the above  
6801       problem.
- 6802   **B.24** Use Plotly in R or Python to plot four June 2015 potential temperature maps  
6803       at the depth layers 5, 25, 105 and 195 meters based on the GODAS data. Try  
6804       to place the four maps together to make a 3D visualization.
- 6805   **B.25** Plot three cross-sectional maps of the December 2015 potential temperature  
6806       at 10°S, equator, and 10°N using the same netCDF GODAS data file as the  
6807       above problem. Each map is on a  $40 \times 360$  grid, with 40 depth levels and 1-deg  
6808       longitude resolution for the equator.
- 6809   **B.26** Do the same as the above but for the June 2015 GODAS potential temper-  
6810       ature data.
- 6811   **B.27** Use R or Python a cross-sectional map of the December 2015 potential tem-  
6812       perature along a meridional line at longitude at 170°E using the same netCDF  
6813       GODAS data file as the above problem.
- 6814   **B.28** Use R or Python to plot a cross-sectional map of the June 2015 potential  
6815       temperature along a meridional line at longitude at 170°E using the same  
6816       netCDF GODAS data file as the above problem.



# Index