

# Contents

## 1 Introduction 1

### 1.1 From data to knowledge: the aim of bioinformatics 1

### 1.2 Using this book 2

- 1.2.1 About the coverage of this book 2
- 1.2.2 Choice of tools 3
- 1.2.3 Choice of operating system 3
- 1.2.4 [www.bixsolutions.net](http://www.bixsolutions.net) 4

### 1.3 Principal applications of bioinformatics 4

- 1.3.1 Sequence analysis 5
- 1.3.2 Transcriptomics 5
- 1.3.3 Proteomics 6
- 1.3.4 Metabolomics 7
- 1.3.5 Systems biology 7
- 1.3.6 Literature mining 8
- 1.3.7 Structural biology 8

### 1.4 Building bioinformatics solutions 8

### 1.5 Publicly available bioinformatics resources 10

- 1.5.1 Publicly available data 10
- 1.5.2 Publicly available analysis tools 14
- 1.5.3 Publicly available workflow solutions 15

### 1.6 Some computing practicalities 16

- 1.6.1 Hardware requirements 16
- 1.6.2 The command line 17
- 1.6.3 Case sensitivity 18
- 1.6.4 Security, firewalls, and administration rights 18

### References 19

## 2 Building biological databases with SQL 21

### 2.1 Common database types 22

- 2.1.1 Flat text files 22
- 2.1.2 XML 23
- 2.1.3 Relational databases 26

### 2.2 Relational database design—the ‘natural’ approach 29

- 2.2.1 Steps 1–3: gather, group, and name the data 30
- 2.2.2 Step 4: data types 35
- 2.2.3 Step 5: atomicity of data 39

when	2.2.4 Steps 6 and 7: indexing and linking tables	39
easier	2.2.5 Departure from design	45
<b>2.3 Installing and configuring a MySQL server</b>	45	
2.3.1 Download and installation	45	
2.3.2 Creating a database and a user account	48	
<b>2.4 Alternatives to MySQL</b>	49	
2.4.1 PostgreSQL	49	
2.4.2 Oracle	50	
2.4.3 MariaDB	50	
2.4.4 Microsoft Access	50	
2.4.5 Big Data and NoSQL databases	51	
<b>2.5 Database access using SQL</b>	52	
2.5.1 Compatibility between RDBMSs	53	
2.5.2 Error messages	53	
2.5.3 Creating a database	53	
2.5.4 Creating tables and enforcing referential integrity	54	
2.5.5 Populating the database	57	
2.5.6 Removing data and tables from the database	59	
2.5.7 Creating and using source files	60	
2.5.8 Querying the database	61	
2.5.9 Transaction handling	68	
2.5.10 Copying, moving, and backing up a database	69	
<b>2.6 MySQL Workbench: an alternative to the command line</b>	70	
<b>2.7 Summary</b>	72	
<b>References</b>	72	

<b>3 Beginning programming in Perl</b>	73
<b>3.1 Downloading and installing Perl</b>	74
3.1.1 Older versions of Perl on Mac OS	74
3.1.2 Older versions of Perl on Linux	75
3.1.3 Installing Perl on Windows	75
3.1.4 Compilers and other developer tools	75
3.1.5 Before getting started	76
<b>3.2 Basic Perl syntax and logic</b>	77
3.2.1 Scalar variables	79
3.2.2 Arrays	85
3.2.3 Hashes	89
3.2.4 Control structures and logic operators	91
3.2.5 Writing interactive programs—I/O basics	97
3.2.6 Some good coding practice	101
3.2.7 Summary	103
<b>3.3 References</b>	103
3.3.1 Multidimensional arrays	104
3.3.2 Multidimensional hashes	107
3.3.3 Viewing data structures with Data::Dumper	110

<b>3.4 Subroutines and modules</b>	112
3.4.1 Making a Perl module	115
<b>3.5 Regular expressions</b>	117
3.5.1 Defining regular expressions	117
3.5.2 More advanced regular expressions	119
3.5.3 Regular expressions in practice	121
<b>3.6 File handling and directory operations</b>	123
3.6.1 Reading text files	124
3.6.2 Writing text files	125
3.6.3 Directory operations	126
<b>3.7 Error handling</b>	127
<b>3.8 Retrieving files from the Internet</b>	129
3.8.1 Utilizing NCBI's eUtilities	131
<b>3.9 Accessing relational databases using Perl DBI</b>	133
3.9.1 Installing DBD::MySQL	134
3.9.2 Connecting to a database	135
3.9.3 Querying the database	136
3.9.4 Populating the database	138
3.9.5 Database transactions and error handling	139
<b>3.10 Harnessing existing tools</b>	140
3.10.1 CPAN	141
3.10.2 BioPerl	142
3.10.3 System commands	143
<b>3.11 Object-oriented programming</b>	143
3.11.1 Object-oriented programming in Perl using Moose	145
<b>3.12 Summary</b>	155
<b>References</b>	156
<b>4 Analysis and visualisation of data using R</b>	157
<b>4.1 Introduction to R</b>	158
4.1.1 Downloading and installing R	159
4.1.2 Basic R concepts and syntax	160
4.1.3 Vectors and data frames	162
4.1.4 The nature of experimental data	165
4.1.5 R modes, objects, lists, classes, and methods	169
4.1.6 Importing data into R	173
4.1.7 Data visualization in R	174
4.1.8 Writing programs in R	180
4.1.9 Some essential R functions	185
4.1.10 The RStudio integrated development environment	189
<b>4.2 Multivariate data analysis</b>	191
4.2.1 Exploratory data analysis	191
4.2.2 Scatter plots	191
4.2.3 Principal components analysis	192

4.2.4 Hierarchical cluster analysis	194
4.2.5 Pattern recognition	198
<b>4.3 R packages</b>	<b>198</b>
4.3.1 Installing and using Bioconductor packages	200
4.3.2 The RMySQL package for database connectivity	205
4.3.3 Packages for multivariate classification	207
4.3.4 Writing your own R packages	207
<b>4.4 Integrating Perl and R</b>	<b>208</b>
<b>4.5 Alternatives to R</b>	<b>208</b>
4.5.1 S+	208
4.5.2 Matlab	209
4.5.3 Octave	210
<b>4.6 Summary</b>	<b>211</b>
<b>References</b>	<b>211</b>

## 5 Developing web resources 213

<b>5.1 Web servers</b>	<b>213</b>
<b>5.2 Introduction to HTML</b>	<b>213</b>
5.2.1 Creating and editing HTML documents	214
5.2.2 The structure of a web page	214
5.2.3 HTML tags and general formatting	215
5.2.4 An example web page	218
5.2.5 Web standards and browser compatibility	220
<b>5.3 Programming for the web using Perl</b>	<b>220</b>
5.3.1 Mojolicious::Lite	221
5.3.2 Debugging Mojolicious applications	224
5.3.3 Routes	225
5.3.4 Interfacing with databases within a web application	227
5.3.5 Getting user input via forms	231
5.3.6 Deploying a Mojolicious application	238
5.3.7 Going further with Mojolicious	239
<b>5.4 Advanced web techniques and languages</b>	<b>239</b>
5.4.1 Cascading stylesheets	239
5.4.2 JavaScript, JavaScript libraries, and Ajax	242
<b>5.5 Data Visualization on the web</b>	<b>244</b>
5.5.1 Using R graphics in Perl	244
5.5.2 Plotting graphs with Chart::Clicker	250
5.5.3 Plotting graphs with SVG::TT::Graph	256
5.5.4 Primitive graphics with Perl	263
5.5.5 Drawing graphs and graphics using JavaScript	263
<b>5.6 Summary</b>	<b>264</b>
<b>References</b>	<b>264</b>

## 6 Software engineering for bioinformatics 265

<b>6.1 Unit testing</b>	<b>266</b>
6.1.1 Unit testing in practice	267
<b>6.2 Version control</b>	<b>272</b>
6.2.1 The basics of version control	272
6.2.2 Centralized versus distributed version control	275
6.2.3 Git	276
6.2.4 Alternatives to Git	286
6.2.5 Hosting and sharing your code on the Internet	287
6.2.6 Running your own code repository	288
<b>6.3 Creating useful documentation</b>	<b>288</b>
6.3.1 Documenting command-line applications	289
6.3.2 Documenting Perl code	290
<b>6.4 User-centred software design</b>	<b>293</b>
<b>6.5 Alternatives to Perl</b>	<b>294</b>
6.5.1 Python	294
6.5.2 Ruby	305
6.5.3 Java	318
6.5.4 Using Galaxy	326
<b>6.6 Summary</b>	<b>327</b>
<b>References</b>	<b>327</b>

## Appendix A: Using command-line interfaces 329

<b>A.1 Getting to the operating system command line</b>	<b>329</b>
<b>A.2 General command-line concepts</b>	<b>331</b>
<b>A.3 Command-line tips</b>	<b>333</b>

## Appendix B: Getting started with Apache HTTP Server 330

<b>B.1 Installing Apache</b>	<b>336</b>
<b>B.2 Apache fundamentals</b>	<b>337</b>

## Appendix C: Setting up a Linux virtual machine in Windows 341

<b>C.1 Installing VirtualBox and configuring a virtual machine</b>	<b>341</b>
<b>C.2 Using the VM</b>	<b>344</b>
<b>C.3 Other uses of virtual machines</b>	<b>345</b>

<b>Index</b>	<b>347</b>
--------------	------------