

Quantitative Bootcamp Tutorials T1-T8

Stephanie Palmer and Stefano Allesina

September 1, 2015

Contents

1	Introduction to R – Basics	1
1.1	What is R?	1
1.2	Why use R?	1
1.3	Launching R	2
1.4	Finding help	2
1.5	R as a calculator	2
1.6	Assignment and data types	4
1.7	Data structures	4
1.7.1	Vectors, matrices, and arrays	4
1.7.2	Lists	9
1.7.3	Data frames	10
1.8	Reading and writing data	11
2	Probability theory	13
2.1	Randomness in biology	13
2.2	Testing your intuition: the bent coin lottery	13
2.3	Binomial distribution	14

2.4	Poisson distribution	15
2.4.1	Interval between events	15
2.5	Gaussian distribution	16
2.6	Winning the bent coin lottery	17
2.7	Entropy	19
2.7.1	Uniqueness	19
2.8	Generating samples of a stochastic process	20
2.9	Markov processes	20
3	Introduction to R – Code flow	21
3.1	The flow of the program	21
3.2	Why writing a script?	21
3.3	Getting started	22
3.4	Branching	22
3.5	Loops	23
3.5.1	Mammals body mass	25
4	Inference	35
4.1	What is inference?	35
4.2	A frightening diagnosis	35
4.3	Bayesian versus Frequentist views of the world and the data in it	36
4.3.1	Priors and posteriors	36
4.3.2	Bayesian view of data	37
4.3.3	Frequentist view of data	37
4.4	Two-envelope paradox	37
4.5	Lindley’s paradox	39
4.6	Regression	40
4.7	Maximum Likelihood (ML) inference	40
4.8	Maximum A Posteriori (MAP) inference	40
4.9	Bayes estimator	41
4.10	Further reading	41
5	Introduction to R – Advanced	43
5.1	Functions	43
5.2	Random numbers	46
5.3	Avoiding loops	46
5.4	Importing code	48
5.5	Importing libraries	48
5.6	Basic plotting	48
5.6.1	Scatterplots: <code>plot</code>	48
5.6.2	Histograms: <code>hist</code>	49

5.6.3	Barplots	50
5.6.4	Boxplots	50
5.6.5	3D plotting (in 2D)	50
5.7	Want to learn more?	51
6	Dynamical systems	53
6.1	What is a dynamical system?	53
6.2	Basic notation	53
6.3	A simple pendulum is not so simple	54
6.4	Phase portraits and stability analysis	54
6.5	Exponential Growth	55
6.5.1	Analysis	56
6.6	The Logistic Growth Model	57
6.6.1	Intraspecific Competition	57
6.6.2	Analysis	57
6.7	Euler method for numerical integration	62
6.8	Further reading	62
7	Introduction to UNIX	63
7.1	What is UNIX?	63
7.2	Why use UNIX?	63
7.3	Directory structure	64
7.4	Using the terminal	65
7.5	Basic UNIX commands	66
7.5.1	How to get help in UNIX	66
7.5.2	Navigating the directory system	66
7.5.3	Handling directories and files	67
7.5.4	Printing and modulating files	67
7.5.5	Time and date	68
7.5.6	Miscellaneous	68
7.6	Advanced UNIX commands	68
7.6.1	Redirection and pipes	68
7.6.2	Selecting columns using <code>cut</code>	70
7.6.3	Substituting characters using <code>tr</code>	71
7.6.4	Selecting lines using <code>grep</code>	73
7.7	Basic scripting	75
8	Data Visualization	79
8.1	Overview	79
8.2	Basic guidelines for displaying data	79
8.3	An example dataset and published plot	80

iv Contents

8.4	A brief introduction to <code>ggplot</code>	80
-----	---	----

Tutorial 1

Introduction to R – Basics

1.1 What is R?

R is a software for statistical analysis. It comes with many built-in functions and excellent graphical capabilities. The main strength of R is that it is fully programmable: you can write code in R and have the software execute it. This means that it is very easy to automate your statistical and data analysis.

The fact that R is easy to program led to the development of thousands of packages, so that you can find a ready-made, specific package for almost any analysis you might want to perform. Because of this strength, R has become the most popular statistical software among biologists.

The main hurdle new users face when learning R is that it is based on a command-line interface: to make things happen, you write text commands in a “shell”, and then the program executes them. This might seem unusual if you come from software based on a Graphical User Interface, where you tend to work by clicking on windows and buttons. However, the command-line is what makes it easy to automate your analysis — all you have to do is collect all the commands in a text file, and then run them in R.

For this brief introduction to R, we are going to use **RStudio**, a graphical interface that simplifies the use of R by giving you immediate access to the code, the shell, and the graphics.

1.2 Why use R?

Writing scripts for all your work, instead of manually typing commands and clicking buttons, makes your research easy to reproduce (just share the scripts with the interested scientists), well-documented (especially if you write meaningful comments to detail what you are doing), and easy to automate (once you have written a script to analyze a dataset, it is easy to make it analyze millions of similar datasets).

2 Tutorial 1 Introduction to R – Basics

R is free software: it is free to use, but it also gives you the freedom to see the code (open source), modify it, and extend it.

1.3 Launching R

Either click on the RStudio icon, or open a terminal and type `rstudio`.

1.4 Finding help

Each command in R comes with a manual page. To access it, type `?NAMEOFCOMMAND` in the console (e.g., `?lm`).

1.5 R as a calculator

To start, we are going to explore some features of R by typing commands in the console. In the console, a “greater sign” (`>`) means that R is ready to accept a command. You can navigate the history of previously typed commands by using the arrows on your keyboard.

Go to the console and type:

```
> 1 + 1
[1] 2
> 1 * 3
[1] 3
> 1.7 * 2
[1] 3.4
> 12 / 3
[1] 4
> 12 / 5
[1] 2.4
> 123 - 72
[1] 51
> 2.1 ^ 5
[1] 40.84101
> log(10)
[1] 2.302585
> log10(10)
[1] 1
> sqrt(9)
[1] 3
> trunc(12.11)
[1] 12
> floor(12.11)
[1] 12
> floor(12.71)
```

```
[1] 12
> trunc(12.71)
[1] 12
> ceiling(12.71)
[1] 13
> round(12.71, 1)
[1] 12.7
```

You can use R as a calculator, with these operators:

Operator	Description
+	addition
-	subtraction
*	multiplication
/	division
^ or **	exponentiation
x %% y	modulus (remainder of integer division)
x %/% y	integer division

R has many built-in mathematical functions:

Function	Description
abs(x)	absolute value
sqrt(x)	square root
ceiling(x)	nearest integer $> x$
floor(x)	nearest integer $< x$
trunc(x)	integer part
round(x, digits=n)	round the number to n digits
cos(x), sin(x), tan(x), etc.	trigonometric functions
log(x)	natural logarithm
log10(x)	base 10 logarithm
exp(x)	e^x

and it can deal with logical values:

```
> 5 > 3
[1] TRUE
> 5 == (10 / 2)
[1] TRUE
> 6 > 2^4
[1] FALSE
> 6 >= (2 * 3)
[1] TRUE
> (5 > 3) & (7 < 5)
[1] FALSE
> (5 > 3) | (7 < 5)
```

```
[1] TRUE
```

1.6 Assignment and data types

When programming in R, you assign values to variables: a variable is a “box” which can contain a variable.

```
> x <- 5
> x * 2
[1] 10
> x <- 7
> x * 2
[1] 14
```

We assigned to the variable `x` the value 5, using the assignment command `<-`. Now we can use `x` to perform operations. If we assign a new value to `x`, the previous value is overwritten.

To list all the variables that you created, type `ls()`. To remove a variable you created, type `rm(NAMEOFVARIABLE)` (e.g., `rm(x)`).

R can handle different types of data:

Type	Description	Example
integer	natural numbers	<code>x <- as.integer(5)</code>
numeric	real numbers	<code>x <- pi</code>
complex	complex numbers	<code>x <- 1 + 3i</code>
logical	TRUE/FALSE	<code>x <- (5 > 7)</code>
character	strings	<code>x <- "hello"</code>

To determine the type of a variable, use the command `class(x)`. You can also test whether a certain variable is of a certain type by using the functions `is.numeric(x)`, `is.character(x)`, etc.

1.7 Data structures

R has several “data structures”, which can be used to organize your data. Each data structure comes with specific operations you can perform.

1.7.1 Vectors, matrices, and arrays

Vectors. The most basic data structure in R is the vector, which is an ordered collection of values. Vectors are defined by concatenating different values with the command `c()`:

```
> x <- c(2, 3, 5, 27, 31, 13, 17, 19)
> x
```



```
[1]  2  3  5 27 31 13 17 19
```

You can access the elements of a vector by their index: the first element is indexed at 1, the second at 2, etc.:

```
> x[3]
[1] 5
> x[8]
[1] 19
> x[9]
[1] NA
```

You can extract several elements at once (i.e., another vector), using the colon (:) command, or by concatenating the indices:

```
> x[1:3]
[1] 2 3 5
> x[4:7]
[1] 27 31 13 17
> x[c(1,3,5)]
[1] 2 5 31
```

You can determine the length of a vector using the function `length(x)`.

Given that R was born for statistics, there are several statistical functions you can perform on vectors (`#` marks comments):

```
> min(x)
[1] 2
> max(x)
[1] 31
> sum(x) # sum all elements
[1] 117
> prod(x) # multiply all elements
[1] 105436890
> median(x) # median value
[1] 15
> mean(x) # arithmetic mean
[1] 14.625
> var(x) # unbiased sample variance
[1] 119.4107
> mean(x^2) - mean(x)^2 # population variance
[1] 104.4844
> summary(x) # print a summary
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 2.00   4.50   15.00   14.62   21.00   31.00
```

You can generate vectors of sequential numbers using the colon command:

```
> x <- 1:10
> x
[1] 1 2 3 4 5 6 7 8 9 10
```

For more complex sequences, use `seq()`:

```
> seq(from = 1, to = 5, by = 0.5)
[1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0
```

To repeat a value or a sequence several times, use `rep()`:

```
> rep("abc", 3)
[1] "abc" "abc" "abc"
> rep(c(1,2,3), 3)
[1] 1 2 3 1 2 3 1 2 3
```

Exercise 1.1

1. Create a vector containing all the even numbers between 2 and 100 (inclusive) and store it in variable `z`.
2. What is the sum of all the elements of the vector?
3. Is it equal to $51 \cdot 50$?
4. Does `seq(2, 100, by = 2)` produce the same vector as `(1:50) * 2`?
5. Extract all the elements of `z` that are divisible by 12. How many elements match this criterion?
6. What happens if you type `z ^ 2`?

Matrices. A matrix is a two-dimensional table of values. In case of numeric values, you can perform the usual operations on matrices (product, inverse, decomposition, etc.):

```
> A <- matrix(c(1, 2, 3, 4), 2, 2) # values, nrow, ncol
> A
      [,1] [,2]
[1,]    1    3
[2,]    2    4
> A %*% A # matrix product
      [,1] [,2]
[1,]    7   15
[2,]   10   22
> solve(A) # matrix inverse
      [,1] [,2]
[1,]   -2  1.5
[2,]    1 -0.5
```

```

> A %*% solve(A)
      [,1] [,2]
[1,]     1     0
[2,]     0     1
> diag(A) # create a vector of the diagonal elements
[1] 1 4
> B <- matrix(1, 3, 2)
> B
      [,1] [,2]
[1,]     1     1
[2,]     1     1
[3,]     1     1
> B %*% t(B) # transpose
      [,1] [,2] [,3]
[1,]     2     2     2
[2,]     2     2     2
[3,]     2     2     2
> Z <- matrix(1:9, 3, 3)
> Z
      [,1] [,2] [,3]
[1,]     1     4     7
[2,]     2     5     8
[3,]     3     6     9

```

To determine the dimensions of a matrix, use `dim()`:

```

> dim(B)
[1] 3 2
> dim(B)[1]
[1] 3
> dim(B)[2]
[1] 2

```

You can access a particular row/column of a matrix:

```

> Z
      [,1] [,2] [,3]
[1,]     1     4     7
[2,]     2     5     8
[3,]     3     6     9
> Z[1,] # first row
[1] 1 4 7
> Z[,2] # second column
[1] 4 5 6
> Z[1:2, 2:3] # submatrix with coefficients in first two rows and
               # second or third column
      [,1] [,2]

```

```
[1,] 4 7
[2,] 5 8
> Z[c(1,3), c(1,3)]
      [,1] [,2]
[1,] 1 7
[2,] 3 9
```

Some operations use all the elements of the matrix:

```
> sum(Z)
[1] 45
> mean(Z)
[1] 5
```

Arrays. If you need tables with more than two dimensions, use arrays:

```
> M <- array(1:24, c(4, 3, 2))
> M
, , 1
      [,1] [,2] [,3]
[1,] 1 5 9
[2,] 2 6 10
[3,] 3 7 11
[4,] 4 8 12

, , 2
      [,1] [,2] [,3]
[1,] 13 17 21
[2,] 14 18 22
[3,] 15 19 23
[4,] 16 20 24
```

You can still determine the dimensions using:

```
> dim(M)
[1] 4 3 2
```

and access the elements as done for the matrices. One thing you should be paying attention to: R drops dimensions that are not needed. So, if you access a “slice” of a 3-dimensional array:

```
> M[, , 1]
      [,1] [,2] [,3]
[1,] 1 5 9
```

```
[2,]    2    6   10
[3,]    3    7   11
[4,]    4    8   12
```

you obtain a matrix:

```
> dim(M[, ,1])
[1] 4 3
```

Exercise 1.2

For this exercise, we're going to use the data you will explore in the Workshop run by Dr. Leslie Osborne. First, we want to be all in the same directory. Hit CTRL+Shift+H: this will open a window asking to choose the working directory. Please go to `QBio/Tutorials/1-T3-T5-Intro_to_R/Sandbox` and click on “Open”. Now that's our new “working directory”. Then, type:
`load("../ ../ ../Workshops/Osborne/Data/MTneuron.RData")`
 to load the data.

The `load()` function loads data that has been previously saved in R. To save your data, simply type `save(obj1, obj2, file = "filename.RData")`, where `obj1` and `obj2` are the objects (you can have as many as you want) to save in the file `filename.RData`.

Note that you can hit **Tab** to auto-complete names, such as the path of a file. You should **always** use **Tab** to complete your paths/filenames, as in this way you can make sure you use long, expressive file names, and that you are not introducing typos.

If everything went well, your workspace should contain:

```
> ls() # list objects in workspace
[1] "directions" "dirtune"    "RFmap"      "theta"
```

1. Determine the dimension of `directions`, and `RFmap`. Are these vectors, matrices, or arrays?
2. How many coefficients are in `RFmap`?
3. What is the mean value of the coefficients in `RFmap`? What is the median?
4. How many elements of `RFmap` are > 0 ?

1.7.2 Lists

Vectors are good if each element is of the same type (e.g., numbers, strings, etc.). Lists are used when we want to store elements of different types, or more complex objects (e.g., vectors, matrices, other lists!). Each element of the list can be referenced either by its index, or by a label:

```

> mylist <- list(Names = c("a", "b", "c", "d"), Values = c(1, 2, 3))
> mylist
$Names
[1] "a" "b" "c" "d"

$Values
[1] 1 2 3

> mylist[[1]] # access first element using index
[1] "a" "b" "c" "d"
> mylist[[2]] # second element using index
[1] 1 2 3
> mylist$Names # using label
[1] "a" "b" "c" "d"
> mylist[["Names"]] # another way to do this
[1] "a" "b" "c" "d"

```

1.7.3 Data frames

Data frames contain data organized like in a spreadsheet. The columns (typically representing different measurements) can be of different types (e.g., a column could be the date of measurement, another the weight of the individual, or the volume of the cell, or the treatment of the sample), while the rows are typically different samples.

When you read a spreadsheet file in R, it is automatically stored as a data frame. The difference between a matrix and a data frame is that in a matrix all the values are of the same type (e.g., all numeric), while in a data frame they can be of different types.

Because typing a data frame by hand would be tedious, let's use a dataset that is already available in R:

```

> data(trees) # Dataset with girth, height and volume of cherry
  trees
> is.data.frame(trees)
[1] TRUE
> dim(trees)
[1] 31  3
> head(trees)
  Girth Height Volume
1   8.3     70   10.3
2   8.6     65   10.3
3   8.8     63   10.2
4  10.5     72   16.4
5  10.7     81   18.8
6  10.8     83   19.7
> trees$Girth

```

```

[1] 8.3 8.6 8.8 10.5 10.7 10.8 11.0 11.0 11.1 11.2 11.3
[12] 11.4 11.4 11.7 12.0 12.9 12.9 13.3 13.7 13.8 14.0 14.2
[23] 14.5 16.0 16.3 17.3 17.5 17.9 18.0 18.0 20.6
> trees$Height[1:5]
[1] 70 65 63 72 81
> trees[1:3,]
  Girth Height Volume
1   8.3     70  10.3
2   8.6     65  10.3
3   8.8     63  10.2
> trees[1:3,]$Volume
[1] 10.3 10.3 10.2

```

Exercise 1.3

1. What is the average height of the cherry trees?
2. What is the average girth of those that are more than 75 ft tall?
3. What is the maximum height of trees with a volume between 15 and 35 ft³?

1.8 Reading and writing data

Reading data frames is very easy: simply use the command `read.table()`, which takes as argument the file name, and can be customized to define a delimiter (space by default), the presence of a header for the column names, etc.

Let's read a file taken from Dr. John Novembre's workshop (make sure you're in the `Sandbox` directory first!):

```
> ch6 <- read.table("../Data/H938_Euro_chr6.geno", header = TRUE)
```

where `header = TRUE` means that we want to take the first line to be a header containing the column names.

How big is this table?

```
> dim(ch6)
[1] 43141      7
```

we have 7 columns, but more than 40k rows! Let's see the first few:

```
> head(ch6)
  CHR      SNP A1 A2 nA1A1 nA1A2 nA2A2
1   6 rs4959515 A  G      0     17    107
2   6  rs719065 A  G      0     26     98
3   6 rs6596790 C  T      0      4    119
4   6 rs6596796 A  G      0     22    102

```

5	6	rs1535053	G	A	5	39	80
6	6	rs12660307	C	T	0	3	121

and the last few:

```
> tail(ch6)
      CHR      SNP A1 A2 nA1A1 nA1A2 nA2A2
43136   6 rs10946282 C  T      0     16    108
43137   6 rs3734763  C  T     19     56     48
43138   6 rs960744   T  C     32     60     32
43139   6 rs4428484  A  G      1     11    112
43140   6 rs7775031  T  C     26     56     42
43141   6 rs12213906 C  T      1     11    112
```

The data contains the number of homozygotes (nA1A1, nA2A2) and heterozygotes (nA1A2), for a number of SNPs obtained by sequencing European individuals:

CHR The chromosome (6 in this case)

SNP The identifier of the Single Nucleotide Polymorphism

A1 One of the alleles

A2 The other allele

nA1A1 The number of individuals with the particular combination of alleles.

Exercise 1.4

1. How many individuals were sampled? Find the maximum of the sum **nA1A1 + nA1A2 + nA2A2**. Notes: you can access the columns by index (e.g., **ch6[,5]**), or by name (e.g., **ch6\$nA1A1**, or also **ch6["nA1A1"]**).
2. You can try using the function **rowSums** to obtain the same result.
3. For how many SNPs do we have that all sampled individuals are homozygotes (i.e., all **A1A1** or all **A2A2**)?
4. For how many SNPs, are more than 99% of the sampled individuals homozygous?

Tutorial 2

Probability theory

2.1 Randomness in biology

This tutorial highlights an important and pervasive aspect of biological systems: stochasticity. (NB: ‘stochasticity’, ‘variability’, ‘uncertainty’, ‘noise’, and ‘fluctuations’ will all be used interchangeably here, though most of these terms have more technical and specific uses in other contexts.) Many of the variables that we observe in biological recordings fluctuate, sometimes because we cannot control all the states of the external and internal world of the organism, other times because thermal noise and other microscopic factors make the state of the biological system we interrogate inherently noisy. It is useful to model not only a median value for a fluctuating variable, but the full shape of its distribution of values.

For example, if we observe the firing of neurons in the brain to repeats of the same external stimulus, the precise times of spikes will vary between repeats. By fitting the statistics of this noise to models we deepen our understanding of the neural response.

In this tutorial, we will cover some basic concepts in probability theory, ending with some fundamental properties of entropy and information. In the Readings folder for this tutorial, you will find a review article by Tkačik and Bialek on information processing in biology. A forthcoming special issue of the Journal of Statistical Physics is also dedicated to this topic.

To build your intuition about quantifying uncertainty, let’s start with a toy problem I first encountered in David MacKay’s lectures on information theory and inference.

2.2 Testing your intuition: the bent coin lottery

A biased coin is used to generate sequences of digits, 1 for heads, 0 for tails, in a lottery. The coin is tossed 25 times to determine the winning sequence. The probability of heads is 0.1. Tickets for the lottery cost \$1 and the prize is \$10,000,000.

Exercise 2.1

1. You are only allowed to purchase one ticket. Which ticket would you buy?

2. How many tickets would you have to buy to cover every possible outcome?

3. Is this lottery worth playing?

2.3 Binomial distribution

Each flip of a coin like this with probability, p , of heads is an example of a Bernoulli trial, the general term for an experiment with only two output states, success or failure. The number of heads in the sequence of independent coin flips generated by our lottery will follow a binomial distribution.

Exercise 2.2

Write down the probability of observing k heads in n coin flips, if the probability of heads is p .

This is the binomial distribution. Rather than memorize this particular form, remember how to write it down as the product of intuitive terms. If we consider the limit of a very small p , we can relate the binomial distribution to the Poisson distribution.

2.4 Poisson distribution

The Poisson distribution describes the probability of finding k events in a fixed interval if we know the rate of occurrence of these events, λ , in that interval. In terms of the variables we have been working with for the bent coin lottery,

(2.1)

We are going to take the limit where p is very small and n is very large, but their product remains fixed.

Exercise 2.3

Derive an expression for the probability of observing k heads in n tosses in the limit of small p or large n .

We have just written down the Poisson distribution. You will see this used as a model for biological variability again and again explicitly or implicitly. It is important to think about whether or not it is a good model for the system under study each time you come across it or are deciding to use it for your own research.

2.4.1 Interval between events

We can also write down the distribution of intervals between events in a Poisson process. This distribution has an exponential form.

Exercise 2.4

Derive an expression for the distribution of an interval, τ , between two events in a Poisson process with rate, λ .

In the limit of large λ , the Poisson distribution is well-approximated by a Gaussian distribution.

2.5 Gaussian distribution

A Gaussian or ‘normal’ distribution (also called a ‘bell-curve’) of a variable x with mean μ and variance σ takes the form

(2.2)

Exercise 2.5

1. As λ gets very large, show that the Poisson distribution approaches a Gaussian distribution with mean λ and variance λ .

You have now derived the form of a very useful distribution and have learned how to do a Gaussian integral along the way.

Exercise 2.6

1. Derive the mean of a Poisson distribution with rate λ :
2. Derive the variance of a Poisson distribution with rate λ :

These quantities are often summarized as the ratio of the mean and variance, or Fano Factor (FF). The FF for a Poisson process is clearly equal to one.

Exercise 2.7

1. Does observing an FF= 1 in data mean that the underlying stochastic process is a Poisson process?

2.6 Winning the bent coin lottery

Now that we have all of these distributions at our fingertips, let's return to the bent coin lottery.

Exercise 2.8

Derive how many tickets you need to buy to guarantee yourself a 99% chance of winning.

One of the new concepts we introduced in this derivation is a particularly fine quantification of uncertainty: entropy.

2.7 Entropy

Measuring uncertainty in a system allows you to attach a single quantity that describes the array of states you observe as its output. This concept of entropy derives its name from analogies to the Boltzmann entropy in statistical physics. Claude Shannon coined this term in his classic paper from 1948 on information theory. Shannon was working on communication systems at the time and was interested in developing a theory that would describe how to capture everything a transmitter produced and the maximal amount a receiver could reliably reproduce, given some noisy channel in between. Information theory has found many uses in biology. In this section, we aim to demystify Shannon's formula for entropy and give you some conceptual tools that will help you critically assess articles that use information theory in biology.

Shannon was interested in a measure, H , of a set of possible events with probabilities p_1, p_2, \dots, p_i , that would quantify uncertainty. He determined that it should obey a few simple rules to be considered a good measure. Namely, from Shannon's original paper:

- **continuity** H should be continuous in p_i .
- **monotonicity** If all p_i are equal, i.e. the distribution over n total states is uniform, such that $p_i = \frac{1}{n}$, H should increase monotonically with n .
- **branching** If a choice is broken down into two successive choices, the original H should be the weighted sum of the individual values of H .

2.7.1 Uniqueness

In this space, write out our derivation, using these three axioms, of the uniqueness of the entropy function, $H = -\sum_i p_i \log p_i$.

2.8 Generating samples of a stochastic process

When modeling biological systems, it is often necessary to generate sequences from a Poisson or other stochastic process. We did this to generate our draws from the bent coin lottery. An introduction to simulating stochastic processes can be found in the Readings folder for this tutorial.

2.9 Markov processes

One feature of the stochastic processes we have been considering today is that they are independent. A flip of the coin doesn't depend on the flip before, or any of the other previous flips. In biological systems, what came before often influences a fluctuating quantity. For example, having spiked, a neuron is unable to spike for a millisecond or two. Modeling this type of variability falls requires using stochastic processes that have an explicit history dependence. Markov processes depend only on the previous time step, in generating the current state. Part of the introduction to point processes in the Readings folder covers Markov processes.

Tutorial 3

Introduction to R – Code flow

3.1 The flow of the program

Now that we're more familiar with R, we turn to writing programs. Typically, you will write your programs in a text file (called a script), with extension `.R`. Then, you can run the script in R by invoking `source(MyScript.R)`.

When you execute the file, R reads the lines of code in order from the top to the bottom. Every time R encounters a command, it will execute it. So in its simplest form, an R program is simply a sequence of commands.

However, it is often important to modify this simple flow of the code: you might have commands that need to be run only if certain conditions are met; commands that need to be run over several files/datasets; commands that you repeat several times; etc.

In this second tutorial on R, we will see how you can modify the flow of a program to suit your needs, and how to organize your code to make it readable and easy to understand.

3.2 Why writing a script?

Before we start, a little motivation on writing scripts. In fact, you can accomplish almost everything you need for your research without writing any — simply type the commands in the shell one at a time. However, organizing your work into well-documented scripts is really important because:

Recycle: you will encounter similar problems in the future, and, having a script, you will be almost done before you even start.

Automate: you will need to repeat the analysis on a different dataset, or slightly tweak it in response to comments. This will take no time at all.

Document: by writing everything in a script, you'll know exactly what you did to get your results. When you'll be writing the Methods section of your paper, you'll be happy you wrote everything down.

Share: believe it or not, people will read what you write, and try to apply your analysis to their own data. Having a script to share will help with this process.

3.3 Getting started

From now on, we'll write scripts, and save them into the **Sandbox** directory within the **T1-T2-T3-Intro.to.R** directory.

In RStudio, hit **CTRL+Shift+N** to start a new script. Organize your scripts well, and save them in the right directory (**Sandbox** in this case), to prevent confusion later on.

To execute the script, click on the **Source** button in the upper-right corner of the area where you see the script, or type `source("Myscript.R")` within the console.

3.4 Branching

The simplest modification of the linear flow of a program is given by conditional branching: if a certain condition is met, then certain commands are executed; otherwise, other commands may be executed.

Let's create a new script called `conditional.R` and save it in the **Sandbox**. Make sure to set the working directory to the **Sandbox**. Now type:

```
1 | z <- readline(prompt = "Enter a number: ")
```

The function `readline()` reads input from the user. It returns a string. Let's convert the string to numeric:

```
2 | z <- readline(prompt = "Enter a number: ")
2 | z <- as.numeric(z)
```

Now we want to determine whether the number is even or odd, and print the answer. If a number z is even, then $z \% 2 == 0$.

```
1 | z <- readline(prompt = "Enter a number: ")
   z <- as.numeric(z)
4 | if (z %% 2 == 0){
   |   print(paste(z, "is even"))
   | } else {
7 |   print(paste(z, "is odd"))
   | }
```

The anatomy of the `if` statement:

```

1 | if (a condition is met){
   |   execute these commands
   | } else {
4 |   execute these other commands [optional]
   | }

```

The `paste` function concatenates strings, and the `print` function prints the results to the console.

Let's try to run the script a few times:

```

> source('conditional.R')
Enter a number: 12
[1] "12 is even"
> source('conditional.R')
Enter a number: 27
[1] "27 is odd"

```

Exercise 3.1

Add code to the script so that:

1. If $z > 100$, the program prints z^3
2. If z is divisible by 17, prints \sqrt{z}
3. If $z < 10$, prints a vector containing the numbers between 1 and z

3.5 Loops

Another way to modify the flow of the program is to write a loop. A loop is simply a series of commands that are repeated a number of times. For example, you want to run the same analysis on all the samples you collected; you want to plot the results contained in a set of files; you want to test your simulation over a number of parameter sets; etc.

R provides you with two ways to loop over commands: the `for` loop, and the `while` loop. Let's start with the `for` loop, which is used to iterate over a vector: for each value of the vector, a series of commands will be run, as shown by the following example, which you can type in a script called `forloop.R`.

```

| myvec <- 1:10 # vector with numbers from 1 to 10
3 | for (i in myvec){
   |   a <- i^2
   |   print(a)
6 | }

```

In the code above, the variable `i` takes the value of each element of `myvec` in sequence. Then, you can use the variable `i` to perform operations.

The anatomy of the `for` statement:

```
for (variable in list_or_vector){
  execute these commands
3 } # automatically moves to the next value
```

For loops are used when you know that you want to perform the analysis using a given set of values (e.g., run over all files of a directory, all samples in your data, all sequences of a fasta file, etc.).

The `while` loop is used when the code is repeated until a certain condition is met, as shown by the following example, which you can type in a script called `whileloop.R`:

```
i <- 1
3 while (i <= 10){
  a <- i^2
  print(a)
6 i <- i + 1
}
```

The script performs exactly the same operations we wrote for the `for` loop above. Note that you need to update the value of `i`, (using `i <- i + 1`), otherwise the loop will run forever (infinite loop — to terminate click on the stop button in the top-right corner of the console).

The anatomy of the `while` statement:

```
while (condition is met){
2   execute these commands
} # beware of infinite loops: remember to update the condition!
```

You can break a loop using the command `break`. For example:

```
i <- 1
3 while (i <= 10){
  if (i > 5){
    break
6  }
  a <- i^2
  print(a)
9 i <- i + 1
}
```

Exercise 3.2

What does this do? Try to guess what each loop does, and then create and run a script to confirm your intuition.

1. Code:

```

1 | z <- seq(1, 1000, by = 3)
2 | for (k in z){
   |   if (k %% 4 == 0){
   |     print(k)
5 |   }
   | }

```

2. Code:

```

1 | z <- readline(prompt = "Enter a number: ")
2 | z <- as.numeric(z)
3 |
   | isthisspecial <- TRUE
   | i <- 1
6 | while (i < z){
   |   if (z %% i == 0){
   |     isthisspecial <- FALSE
9 |     break
   |   }
   | }
12 |
   | if (isthisspecial == TRUE){
   |   print(z)
15 | }

```

3.5.1 Mammals body mass

Now we are going to explore some features of a dataset detailing the body mass of mammals of the late Quaternary (Smith *et al.*, Ecology 2003). The data is contained in the file `MOMv3.3.txt`, which you can find in the `Data` directory of the tutorial, along with the metadata.

Our goal is to calculate the average body mass (in grams) for each of the 152 Families of mammals represented in the data.

We are going to do this in three different ways: the good, the bad, and the ugly. We will get to exercise our ability to write for loops, and see how different designs of the program (all correct), have different time requirements.

Before starting, let's do the operations we will need to repeat for each approach. Reading the data:

```
> dd <- read.table("../Data/MOMv3.3.txt", header = FALSE, sep = "\t",
  stringsAsFactors = FALSE)
> dim(dd)
[1] 5731    9
```

The special command `stringsAsFactors = FALSE` tells R that we do not want to treat the strings in the data as **factors** (levels, very useful for linear regressions, but not for our purposes). Next, we assign a header (the Authors did not provide one):

```
> colnames(dd) <- c("Continent", "Status", "Order", "Family", "Genus",
  "Species", "LogMass", "CombinedMass", "Reference")
> head(dd)
```

	Continent	Status	Order	Family	Genus	Species
1	AF	extant	Artiodactyla	Bovidae	Addax	nasomaculatus
	4.85	70000.3	60			
2	AF	extant	Artiodactyla	Bovidae	Aepyceros	melampus
	4.72	52500.1	63, 70			
3	AF	extant	Artiodactyla	Bovidae	Alcelaphus	buselaphus
	5.23	171001.5	63, 70			
4	AF	extant	Artiodactyla	Bovidae	Ammodorcas	clarkei
	4.45	28049.8	60			
5	AF	extant	Artiodactyla	Bovidae	Ammotragus	lervia
	4.68	48000.0	75			
6	AF	extant	Artiodactyla	Bovidae	Antidorcas	marsupialis
	4.59	39049.9	60			

Column 7 contains the log body mass (in grams) of an adult. Note that the authors used the value `-999` to denote the missing data. In R, it is better to use `NA` (Not Available), as there are special methods for dealing with missing data.

```
> dd[dd == -999] <- NA
> summary(dd[,7])
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.260	1.360	2.070	2.484	3.410	8.280	1372

Finally, let's extract the names of all the families in the dataset:

```
> Fam <- sort(unique(dd$Family))
> Fam
[1] "Abrocomidae" "Acrobatidae" "Agoutidae"
"Anomaluridae"
...
> nfam <- length(Fam)
```

```
> nfam
[1] 152
```

where the function `unique` returns all the unique elements of a vector.

Ugly

We start with our “ugly” (albeit perfectly correct!) strategy. First, copy and paste these commands in a .R script called `ugly.R` (to see all recent commands, type `history()`):

```
1 # read the data
  dd <- read.table("../Data/MOMv3.3.txt", header = FALSE,
                    sep = "\t", stringsAsFactors = FALSE)
4 # assign column names for header
  colnames(dd) <- c("Continent", "Status", "Order", "Family",
                    "Genus", "Species", "LogMass", "CombinedMass", "
                    Reference")
7
  # set missing data to NA
  dd[dd == -999] <- NA
10
  # store number of rows
  nrows <- dim(dd)[1]
13
  # extract unique families
  Fam <- sort(unique(dd$Family))
16
  # store number of families
  nfam <- length(Fam)
```

Now, let's outline the ugly strategy:

```
...
3 # store number of families
  nfam <- length(Fam)

6 # 1) create an empty data frame
  # 2) for each family
  # 3) cycle through all the records:
9 # - if the species belongs to the family
  # - if the body mass is not NA
  # - add to the avg mass (after exponentiating)
```

```
12 # – keep track of the number of species in the family
    # 4) once finished, average and store the result
```

Let's start writing the code:

```
...
2 # 1) create an empty data frame
  avg_BodyMass_Family <- data.frame()
...
```

and let's start working on the first for loop:

```
    # 2) for each family
2  for (f in 1:nfam){
    my_family <- Fam[f]
    my_avg_bodymass <- 0
5    my_numberofspecies <- 0
    print(my_family)
    # 3) cycle through all the records:
8    # – if the species belongs to the family
    # – if the body mass is not NA
    # – add to the avg mass (after exponentiating)
11   # – keep track of the number of species in the family
    # 4) once finished, average and store the result
  }
```

Now run the code to see that everything is good so far:

```
> source('ugly.R')
[1] "Abrocomidae"
[1] "Acrobatidae"
...
```

We can finish the code:

```
    # read the data
2  dd <- read.table("../Data/MOMv3.3.txt", header = FALSE,
                   sep = "\t", stringsAsFactors = FALSE)
    # assign column names for header
5  colnames(dd) <- c("Continent", "Status", "Order", "Family",
                   "Genus", "Species", "LogMass", "CombinedMass", "
                   Reference")

8  # set missing data to NA
  dd[dd == -999] <- NA
```


[illegible]

```

53     }                                     NumSpecies = my_numberofspecies))
    }
    # print the results
56 print(avg_bodymass_family)

```

If your code is correct, you should obtain:

	Family	AvgBodyMass	NumSpecies
1	Abrocomidae	2.202889e+02	3
2	Acrobatidae	4.365158e+01	1
3	Agoutidae	8.266358e+03	3
4	Anomaluridae	4.848419e+02	7
...			

The code is very slow, as measured by running:

```

> system.time(source("ugly.R"))
...
   user  system elapsed
307.630   0.000  307.815

```

On the computer I am using to write these notes, it takes more than 5 minutes! This is not surprising, given that for each family (152), we go through all the records in the database (5731), and therefore the computer is doing more than 800,000 operations!

In the next section, we modify the code to reduce the running time.

Bad

Copy the beginning of the `ugly.R` file into `bad.R` file, and let's outline the strategy:

```

# read the data
2 dd <- read.table("../Data/MOMv3.3.txt", header = FALSE,
                  sep = "\t", stringsAsFactors = FALSE)
# assign column names for header
5 colnames(dd) <- c("Continent", "Status", "Order", "Family",
                  "Genus", "Species", "LogMass", "CombinedMass", "
                  Reference")

8 # set missing data to NA
dd[dd == -999] <- NA

11 # store number of rows
nrows <- dim(dd)[1]

```

```

14 # extract unique families
    Fam <- sort(unique(dd$Family))

17 # store number of families
    nfam <- length(Fam)

20 # 1) create a dataframe with the names of each
    # family, and initialize AvgBodyMass and NumSpecies
    # to 0 for each row.

23 # 2) cycle only once through the data, and,
    # for each species, update values for
26 # the corresponding family

    # 3) remove the families with 0 species
29 # 4) average the results

```

Now let's fill in the code:

```

# 1) create a dataframe with the names of each
# family, and initialize AvgBodyMass and NumSpecies
3 # to 0 for each row.
    avg_bodymass_family <- data.frame(Family = Fam,
                                     AvgBodyMass = 0,
6                                     NumSpecies = 0)
    print(avg_bodymass_family)

```

and run it to make sure everything is good.

Next, we write the cycle:

```

# 2) cycle only once through the data, and,
2 # for each species, update values for
    # the corresponding family
    for (i in 1:dim(dd)[1]){
5      my_family <- dd[i,]$Family
        if (is.na(my_family) == FALSE){
            # find which row we need to update
8      my_row <- which(Fam == my_family)
            # if we have a body mass
            if (is.na(dd[i,]$LogMass) == FALSE){
11             # update the values

```

```

    avg_bodymass_family[my_row, ]$NumSpecies <- avg_bodymass_family[my_
      row, ]$NumSpecies + 1
    avg_bodymass_family[my_row, ]$AvgBodyMass <- avg_bodymass_family[my_
      _row, ]$AvgBodyMass + 10 ^ dd[i,]$LogMass
14   }
    }
  }
17 print(avg_bodymass_family)

```

Finally, remove the families with zero species, and average:

```

1 # 3) remove the families with 0 species
avg_bodymass_family <- avg_bodymass_family[avg_bodymass_family$
  NumSpecies > 0, ]
# 4) average the results
4 avg_bodymass_family$AvgBodyMass <- avg_bodymass_family$AvgBodyMass / avg
  _bodymass_family$NumSpecies
print(avg_bodymass_family)

```

Running the code, we find that it is much, much faster:

```

> system.time(source("bad.R"))
...
      Family AvgBodyMass NumSpecies
1   Abrocomidae 2.202889e+02         3
2   Acrobatidae 4.365158e+01         1
3   Agoutidae 8.266358e+03         3
...
   user  system elapsed
7.159   0.000   6.471

```

The fact that this is much faster is not surprising, as we are doing about 6,000 operations instead of 800,000. We can do even better, though.

Good

Copy the beginning of the file into the new script `good.R`, and let's outline our strategy:

```

# read the data
dd <- read.table("../Data/MOMv3.3.txt", header = FALSE,
3   sep = "\t", stringsAsFactors = FALSE)
# assign column names for header
colnames(dd) <- c("Continent", "Status", "Order", "Family",
6   "Genus", "Species", "LogMass", "CombinedMass", "
    Reference")

```

```

# set missing data to NA
9 dd[dd == -999] <- NA

# store number of rows
12 nrows <- dim(dd)[1]

# extract unique families
15 Fam <- sort(unique(dd$Family))

# store number of families
18 nfam <- length(Fam)

# 1) create an empty dataframe
21 # 2) for each family
# - find the subset of LogMass for that family
# - add a record to the dataframe
24 # 3) print the results

```

Let's fill in the code: we're going to use the function `subset` to extract the records that match our criteria

```

# 1) create an empty dataframe
avg_bodymass_family <- data.frame()
3

# 2) for each family
for (f in Fam){
6 # - find the subset of LogMass for that family
ddsub <- subset(dd$LogMass, dd$Family == f & is.na(dd$LogMass) ==
FALSE)
# - add a record to the dataframe
9 if (length(ddsub) > 0){
avg_bodymass_family <- rbind(avg_bodymass_family,
data.frame(Family = f,
12 AvgBodyMass = mean(10 ^ ddsb),
NumSpecies = length(ddsub)))
}
15 }

# 3) print the results
print(avg_bodymass_family)

```

Now running the code takes less than a second:

```
> system.time(source("good.R"))
      Family AvgBodyMass NumSpecies
1      Abrocomidae 2.202889e+02      3
2      Acrobatidae 4.365158e+01      1
3      Agoutidae 8.266358e+03      3
...
      user  system elapsed
0.000    0.000    0.287
```

Again, this is not surprising, as our loop involves cycling over 152 values instead of the original 800,000.

What have we learned?

The long example above shows you that there isn't a single, correct way of writing a program. Because we are scientists, we care very much about the correctness of the code — our results only hold if the program we write does exactly, and solely, what is intended to do. Hence, it is good practice to write the same script in multiple, alternative ways, and make sure that all of them return the same result.

Different strategies have different costs and advantages. Time is often a limiting issue: you want to be able to obtain your results in a reasonable amount of time. Remember that the total time it takes to get results is the sum of three elements: a) the time it takes to write the code; b) the time it takes to debug the code (i.e., make sure it is correct); c) the time it takes to run the code. Most novices put too much emphasis on c), while they are spending much longer on a) and b).

Another element you want to take into account when deciding your strategy is the readability of the code. Often, being too clever and condensing complex operations in one command makes the code difficult to read. Hence, it is better to have a code that runs in 5 seconds (which is a drop in the ocean of the time you spend working on a paper), but that is easy to read and understand, than a code that runs in a few milliseconds, but which is so clever that next week you'll not be able to understand what you did.

Exercise 3.3

1. Write a script that stores, for each family, the maximum body size of extinct and extant species.
2. Find the proportion of families for which the largest extinct species is larger than the largest extant species.

Tutorial 4

Inference

4.1 What is inference?

Inference is the act of drawing conclusions from data, usually by making some assumptions about the structure of the data. This involves selecting a model that describes how the data were generated and then drawing some conclusions (inferences) about this model, given the sampled data. Scientists in the machine learning community sometimes use the term ‘inference’ to describe the particular process of finding the time-evolving, unobserved or ‘hidden’ states in their models. More generally, scientists use the term inference to refer to the act of fitting statistical models to data. These can be fully parametric or non-parametric or mixed. Our goal in this tutorial is to familiarize ourselves with Bayesian and frequentist approaches to data, highlighting which approach is more useful given the problem we are trying to solve.

4.2 A frightening diagnosis

Let’s start with a simple problem to build our intuition about how to make inferences from data.

Exercise 4.1

Hester is given a test for a terrible disease. The result of this test can be only positive (indicating presence of the disease) or negative. The test gives accurate positive results for 95% of those tested who have the disease, and accurate negative results for 95% of those tested who do not have the disease. About 0.5% of people in Hester’s demographic have the disease. The test returns a positive result for Hester.

1. What is the probability that poor Hester has the disease?

To combine all of the information given carefully, we will use a simple identity dubbed Bayes' Rule to write out the relevant conditional probability distributions. Follow the derivation supplied at the board and write it down here:

4.3 Bayesian versus Frequentist views of the world and the data in it

4.3.1 Priors and posteriors

In the previous derivation, we used Bayes' Rule to help us construct the correct estimate of Hester's probability of disease. In the denominator, we averaged over something called the prior distribution. The prior incorporated our knowledge of the risk for the disease in Hester's demographic, before we had knowledge of her test result. In this case, the prior greatly modified our estimate of whether or not Hester had the disease. The probability we calculated is called the posterior. It measures the probability of disease given Hester's test result. Moving from a prior to a posterior value by incorporating data is called a Bayes update. Moving from a prior distribution to a posterior distribution is a true Bayesian step.

4.3.2 Bayesian view of data

In a Bayesian framework, the parameters that one is trying to estimate are characterized by probability distributions that one has beliefs about characterized by prior distributions. A Bayesian uses data to answer the question: Which parameters are most likely? The Bayesian view of the world is, in some sense, very abstract. There are no real fixed values of parameters in the world, only distributions. Bayesians must often perform averages over distributions of parameters to arrive at estimates. Because of this, it is sometimes joked that Bayesians spend their lives doing integrals.

4.3.3 Frequentist view of data

In the frequentist view of the world, there are true underlying physical parameters that have specific values, which we are trying to estimate from random samples of data. Frequentists set parameters with the data. Frequentists often ‘average over the data’ while Bayesians ‘average over parameter distributions’.

We will now explore two classic problems that illustrate the differences between these approaches and will build your intuition about when to use each approach.

4.4 Two-envelope paradox

We begin by laying out the problem, which was first described over 50 years ago by Maurice Kraitchik.

Exercise 4.2

Two identical envelopes are prepared. One contains a quantity of money, x , while the other contains twice as much, $2x$. You pick an envelope, but before opening it or otherwise gaining any information about its contents, you are asked if you would like to switch envelopes or keep the one you have.

1. Should we stay or switch?

A simpler scenario, the so-called necktie paradox, may help clarify your thinking about this problem. Two men are at a Father’s Day party and both have been given ties by their children. They argue over which tie was more expensive. They propose a bet. They will consult their children and find out whose is pricier. The one with the more expensive tie has to give it to the other man. Each dad’s reasoning goes

like follows: I have a 50/50 chance of winning. If I bet and lose, I lose the value of my tie. If I bet and win, I gain more than the value of my tie. Paradoxically, both men seem to have an advantage in betting.

2. Can you describe the flaw in this logic?
3. In the two-envelope paradox, describe a similar string of flawed logic:
4. Let's describe the apparent paradox in the two-envelope problem: (copy from the board)
5. What would a Bayesian do?

The two-envelope problem shows us how a knee-jerk approach leads us astray, but a Bayesian approach, explicitly calculating conditional probabilities, resolves the apparent paradox.

4.5 Lindley's paradox

We now turn to another apparent paradox, that will show us how Bayesian and frequentist approaches can arrive at opposing conclusions. Lindley's paradox is about model comparison between H_0 , our null, and H_1 , given some data, x . It exists when a frequentist rejects the null hypothesis, H_0 , but a Bayesian favors H_0 over H_1 .

Exercise 4.3

A classic example of Lindley's paradox applied to estimating the boy/girl birth ratio in a population. In the city Bayfreak, 49,581 boys and 48,870 girls were born in the last three years. Assume that the number of male births is a binomial variable with parameter, θ . We wish to test whether $\theta = 0.5$ or some other value.

1. What would a frequentist do?

2. What would a Bayesian do?

Lindley's paradox teaches us that if you have information, use it, but if you do not, don't make an overly diffuse prior. If you do, you will need a lot of data to overcome it.

4.6 Regression

The term regression describes a model class for performing inference. The 'linear' part of linear regression describes the structure of the relationship between the data and the parameters in the model. If your function is linear in the fitted parameters, it is linear regression, even if the model being fit is, say, a polynomial where the parameters are the coefficients in front of each term. Let us label our sampled data, x , and our the parameters we would like to fit, λ . We will now cover some of the most popular methods for estimating λ .

4.7 Maximum Likelihood (ML) inference

In the maximum likelihood framework, we seek to find the λ that maximize

(4.1)

the likelihood that the data, x , were drawn from a distribution defined by the parameters, λ . There is an important distinction between a probability distribution and a likelihood, though they may be written in the same form. Take this conditional probability, $P(x|\lambda)$, as an example. If the λ are fixed numbers, then this is just a probability distribution over x . If the λ are the variables and the data are fixed, then $P(x|\lambda)$ is a likelihood. This likelihood is not a probability distribution over the parameters.

It is often useful to work with the \ln of a function, and maximizing a function or its logarithm are equivalent.

Within the ML framework, one can add a Bayesian prior that the parameters are most likely zero, amounting to a type of L1 penalty.

4.8 Maximum A Posteriori (MAP) inference

In contrast, maximum a posteriori inference seeks to maximize the conditional probability of the parameters given the data:

(4.2)

We use Bayes' Rule to express this quantity in terms of things we can measure. Expressing $P(\lambda|x)$ in this way, we have

(4.3)

The $P(x)$ are the same in both the ML and MAP frameworks. What MAP inference has added to the mix is the prior on the parameters, $P(\lambda)$. This addition might lead you to

infer that MAP inference is Bayesian, however MAP inference is a point estimator (i.e. the output of the procedure is a set of numbers that are the fit parameters) while the output of a Bayesian estimator would be characteristics of a distribution of parameters.

Exercise 4.4

1. Show how ML and MAP are related when the prior is uniform:

4.9 Bayes estimator

A Bayesian estimator seeks to minimize the ‘risk’ generated by a loss function, $L(\lambda, \lambda^{\text{est}})$, in forming an estimate of the parameters, λ^{est} , given that the true parameters are λ . This risk is

(4.4)

where x is, again, the data. The mean-squared error is a commonly used loss function.

4.10 Further reading

In this tutorial, we have focused on big concepts rather than particular methods for model analysis, generation, and data sampling. Some tutorials on the techniques that you should familiarize yourself with are included in the Readings section of this tutorial and cover: Hidden Markov models, ROC analysis, quantile-quantile or QQ plots, and the Markov chain Monte Carlo (MCMC) sampling technique.

Tutorial 5

Introduction to R – Advanced

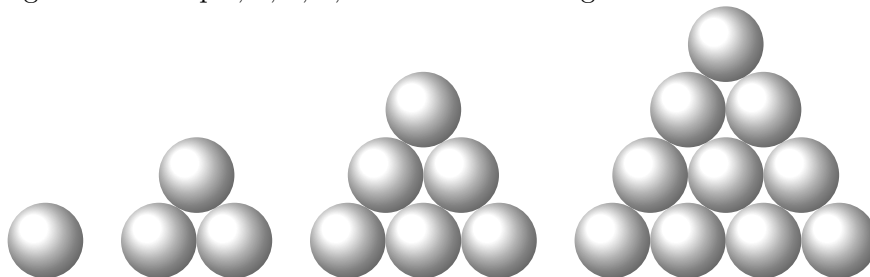
5.1 Functions

In the previous tutorials, we used many built-in functions (e.g., `which`, `length`, `dim`). R allows you to write your own functions, and call them within your programs.

The anatomy of a function is:

```
1 my_function_name <- function([optional arguments]){  
  operations  
  return(value_to_return) [optional]  
4 }
```

To start working with functions, we will write one that tells us whether a number is “triangular”. Triangular numbers count objects that can be arranged in an equilateral triangle. For example, 1, 3, 6, and 10 are a triangular numbers:



Each triangular number T can be written as $T = n(n + 1)/2$ (e.g., $T = 1, n = 1$; $T = 3, n = 2$). Hence, an integer y is triangular if

$$n = (\sqrt{8y + 1} - 1)/2 \tag{5.1}$$

is also an integer. We can write a function that checks whether an integer is triangular:

```
| is.triangular <- function(y){
```

```

2 | # checks whether a number is triangular
   | # if so, then n
   | n <- (sqrt(8 * y + 1) - 1) / 2
5 | # should be an integer
   | if (as.integer(n) == n) {
   |   return(TRUE)
8 | }
   | return(FALSE)
   | }

```

Save this function in the script `triangular.R`. We can make sure that everything is working:

```

> source("triangular.R") # read the script
> is.triangular(91)
[1] TRUE
> is.triangular(78)
[1] TRUE
> is.triangular(4)
[1] FALSE
> is.triangular(56)
[1] FALSE

```

Now let's write another function that returns all the triangular numbers between 1 and `max.number`:

```

# Now find and store the triangular numbers
find.triangular <- function(max.number){
3 | # strategy:
   | # iterate using a for loop
   | # build result vector by concatenating
6 | to.test <- 1:max.number
   | triangular.numbers <- numeric(0) # initialize empty vector
   | for (i in to.test){
9 |   if (is.triangular(i)){
   |     triangular.numbers <- c(triangular.numbers, i)
   |   }
12 | }
   | print(paste("There are", length(triangular.numbers),
   |           "triangular numbers between 1 and ", max.number))
15 | return(triangular.numbers)
   | }

```

Testing:


```

> source("triangular.R") # read the script
> find.triangular(100)
[1] "There are 13 triangular numbers between 1 and 100"
[1] 1 3 6 10 15 21 28 36 45 55 66 78 91
> find.triangular(1000)
[1] "There are 44 triangular numbers between 1 and 1000"
[1] 1 3 6 10 15 21 28 36 45 55 66 78 ...
> find.triangular(10000)
[1] "There are 140 triangular numbers between 1 and 10000"
[1] 1 3 6 10 15 21 28 36 45 ...

```

Here we wrote two functions each taking an argument (y for `is.triangular`; `max.number` for `find.triangular`). You can also write functions that do not require arguments; if you have several arguments, separate them using commas.

You can return only **one** value. In case you need to return multiple things, put them in a list/vector and return it.

```

# a function with no arguments and no return
2 tell.fortune <- function(){
  if (runif(1) < 0.3) {
    print("Today is going to be a great day for you!")
5  } else {
    print("You should have stayed in bed")
  }
8 }

# a function taking multiple values
11 # and returning a vector
order.three <- function(a, b, c){
  return(sort(c(a, b, c)))
14 }

# a function taking multiple values
17 # and returning a list
order.three.list <- function(a, b, c){
  ordered <- sort(c(a,b,c))
20  return(list("First" = ordered[1],
              "Second" = ordered[2],
              "Third" = ordered[3]))
23 }

```

Exercise 5.1

1. Write a function that takes three arguments, `x1`, `x2`, and `x3`, and determines whether their sum is a pentagonal number. Pentagonal numbers are integers that can be written as $P = n(3n-1)/2$. Thus, the integer y is pentagonal if $x = (\sqrt{24y+1}+1)/6$ is also an integer. For example, 1, 5, 12, 22, and 35 are the first few pentagonal numbers.

5.2 Random numbers

R can sample (pseudo-)random numbers from many distributions. This is very useful for simulations! Examples:

```
> runif(3) # extract three numbers from uniform U[0,1]
[1] 0.8252214 0.8811069 0.8099231
> rnorm(4, mean = 1, sd = 5) # Normally distributed
[1] 7.2729430 0.3416511 5.5701139 8.5061745
> rpois(4, lambda = 5) # Poisson with rate lambda
[1] 4 7 4 7
```

You can easily sample from a vector:

```
> v <- 1:10
> sample(v, 2) # sample without replacement from v
[1] 5 8
> sample(v, 11, replace = TRUE) # sample with replacement
[1] 9 4 8 8 2 4 6 1 7 8 10
```

5.3 Avoiding loops

Sometimes loops in R can be very slow. Because of this, R comes with a number of functions meant to avoid using loops. For example, you can do whole-vector operations:

```
> A <- runif(100000)
> B <- rnorm(100000)
>
> system.time({
+   Z <- numeric(0)
+   for (i in 1:length(A)){
+     Z <- c(Z, A[i] + B[i])
+   }
+ })
   user  system elapsed 
23.589   0.000   23.607 
>
> system.time(Z <- A + B)
```

```

user  system elapsed
0      0      0

```

```

> system.time({
+   my_sum <- 0.0
+   for (i in A){
+     my_sum <- my_sum + i
+   }
+ })
  user  system elapsed
0.035   0.000   0.035
>
> system.time(sum(A))
  user  system elapsed
0      0      0

```

And there are special functions to do row/column operations:

```

> GetRowMeans <- function(M){
+   myRowMeans <- rep(0, dim(M)[1])
+   for (i in 1:dim(M)[1]){
+     myRowMeans <- mean(M[i,])
+   }
+   return(myRowMeans)
+ }
>
> M <- matrix(runif(10000 * 10000), 10000, 10000)
>
> system.time(GetRowMeans(M))
  user  system elapsed
2.292   0.241   2.651
> system.time(rowMeans(M))
  user  system elapsed
0.260   0.008   0.269
> system.time(colSums(M))
  user  system elapsed
0.104   0.005   0.729

```

You can use `apply` to “apply” a function to a matrix/array along a certain dimension:

```

> M <- array(rnorm(5 * 7 * 3), c(5, 7, 3))
> apply(M, 3, sd)
[1] 0.9713697 0.8401228 0.8312546
> apply(M, 1, sum) # equivalent to rowSums
[1] 7.28416218 5.16152143 0.03278329 -1.15950951 -1.89666137
> apply(M, 2, mean) # equivalent to colMeans

```

```
[1] 0.2585351 0.1696397 0.0707347 -0.1895531 0.3822819
     0.2574750 -0.3209603
```

There are methods to apply a certain function to all elements of a list (`lapply`), a vector, etc. If you need to do this type of operations often, it is worth checking the package `plyr` out.

5.4 Importing code

If you are writing functions that you are using a lot, then it is convenient to save them into an `.R` file, and then import them in your other code. You can do this by including

```
| source("file_to_import.R")
```

into your scripts. Then, you can use all the functions contained in `file_to_import.R` as if they were written in your current script.

5.5 Importing libraries

Similarly, you can take advantage of the many packages that are available for R. These can be installed from `RStudio` using the Packages tab you find in the lower-right panel. Once you install a package, you can access all the functions contained in the package by including the line

```
| library(name_of_package)
```

into your scripts.

5.6 Basic plotting

In R, you can choose among different ways of plotting your data. By default, R ships with built-in graphical functions. Other functions, or entirely different paradigms, can be enabled using different packages. Here we give a brief overview of the standard plotting. In Tutorial 8 you will get to work with `ggplot`, which is a very powerful alternative to the standard plots in R.

In standard R graphics, plots are drawn interactively: you can start with a plot, and then overlay other elements on the existing plot.

5.6.1 Scatterplots: `plot`

This is the most basic function to plot points defined by their x and y coordinate.

```
| # Simple plot
2 | x <- 1:30
```

```

y <- rnorm(30, mean = x)
y2 <- rnorm(30, mean = x, sd = sqrt(x))
5 # plot y against x
  plot(y ~ x)
  # alternatively
8 plot(x, y)
  # using lines instead of points
  plot(x, y, type = "l")
11 # using both
  plot(x, y, type = "b")
  # change the type of point
14 plot(x, y, type = "b", pch = 4)
  # change color
  plot(x, y, type = "b", pch = 4, col = "blue")
17 # add a line
  abline(c(0, 1)) # intercept, slope
  # add another data set
20 points(x, y2, col = "red")
  # set x-label and y-label
  plot(x, y2, col = "orange", xlab = "my x label", ylab = "yyy")
23 # set ranges
  plot(x, y2, xlim = c(1,10))

```

5.6.2 Histograms: hist

The function `hist` is used to produce simple histograms:

```

# Histograms
d1 <- rpois(100, lambda = 3)
3 # basic histogram
  hist(d1)
  # specify desired number of bins
6 hist(d1, breaks = 4)
  # specify bin edges
  hist(d1, breaks = c(0, 1, 3, 5, 7, 11, 21))
9 # use frequencies (default if bins have equal size)
  hist(d1, freq = TRUE)
  # use density
12 hist(d1, freq = FALSE)
  # get the histogram, but without plotting it
  z <- hist(d1, plot = FALSE)

```

```

15 | # access elements of the histogram
    | z$counts # counts per bin
    | z$mids  # midpoints of the bins

```

5.6.3 Barplots

Barplots are used to represent data for different groups:

```

1 | data(islands) # area of islands
  | barplot(islands)
  | barplot(islands, horiz = TRUE) # horizontal
4 | barplot(islands, horiz = TRUE, las = 1) # change orientation of labels

  | data(iris)
7 | barplot(height = iris$Petal.Width, beside = TRUE, col = iris$Species)

```

5.6.4 Boxplots

Boxplots are used to show the range of a distribution, and the location of the bulk of its mass:

```

  | data(iris)
2 | boxplot(iris$Petal.Width ~ iris$Species, col = c("red", "green", "blue")
  |         )

```

5.6.5 3D plotting (in 2D)

To show density plots, or plot a matrix:

```

1 | # 3D plotting in 2D
  | x <- sort(rnorm(100))
  | y <- sort(rnorm(50))
4 | z <- x %o% y # outer product
  | image(z) # very simple
  | filled.contour(z) # fancier

```

5.7 Want to learn more?

If so, you're in luck: as soon as you get back to Chicago, you could participate in the **R Software Carpentry** workshop offered to all PhD students, postdocs, and researchers in the BSD.

The workshop is sponsored by the Office of Graduate and Postdoctoral Affairs and the Department of Human Genetics, and is taught by John Blischak (jdblischak@uchicago.edu) and Emily Davenport (ed379@cornell.edu).

When: September 15-16, 2015

Where: Stuart 101, 5835 South Greenwood Avenue Chicago, IL 60637.

What: dplyr, ggplot, Git, debugging, programming style, writing reproducible reports.

How: Go to <https://2015-09-15-chicago.eventbrite.com/> and register (pwd uchicago).

Register right now, as enrollment is limited.

Tutorial 6

Dynamical systems

6.1 What is a dynamical system?

A dynamical system is anything that changes in time. In biology, dynamics are important from the molecular to the ecological scale, and at every level in between. Often, the time evolution of a biological system is complex. This tutorial will introduce you to some basic and powerful tools for analyzing dynamics in nonlinear systems.

6.2 Basic notation

We begin by defining a general dynamical system with a set of n variables, x_i . The dynamics of this collection of variables is generally written as

(6.1)

where $\dot{x} = \frac{dx}{dt}$, and $\vec{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ is called the ‘phase space’ or ‘state space’ of the dynamical system. This equation can be broken out, component-wise, so that we have an equation for each variable

(6.2)

The system is linear if all of the x_i ’s on the right hand side appear to first power only. Some examples of nonlinear terms are: x_1^2 , $\sin x_1$, and $x_1 x_3$. Note that we can transform non-autonomous systems, those with explicit time dependence, into this form via:

(6.3)

6.3 A simple pendulum is not so simple

From our first course in general physics, we learn to solve for the equation of motion of a simple harmonic oscillator, $m\ddot{x} + kx = 0$, arriving at the familiar equation, $x(t) = A \cos(\sqrt{\frac{k}{m}}t + \phi)$. We can convert this nonlinear-looking second-order equation into the above form by making the substitution $x_1 = x$ and $x_2 = \dot{x}$ to arrive at a coupled system of linear equations

(6.4)

Here we see the power of writing out the equations component-wise. In this formulation the two variables now correspond to the position and velocity of the oscillator. These equations only describe the motion of a real pendulum for small perturbations around 0. The implicit approximation made here is that $\sin x \approx x$ for small x . The true equation for a pendulum is $\ddot{x} + \sin x = 0$. We can also write this out component-wise

(6.5)

and we see that this is a nonlinear dynamical system that looks quite difficult or impossible to solve. Oftentimes in biological systems, your equations will look like this and seem quite intractable. What can we do when faced with a problem like this?

6.4 Phase portraits and stability analysis

The big idea here is to run our construction in reverse. Instead of solving for the equations that govern the motion of our variables, we instead try to figure out what trajectories look like in phase space and thereby gain knowledge about those solutions. We will do that all without solving the equations! Instead, we will construct a ‘phase portrait’ of all of the qualitatively different trajectories in the phase space.

As an example, consider a one-dimensional version of the pendulum equation we just encountered, $\dot{x} = \sin x$, and let’s try to solve for $x(t)$ in the usual way and then try to analyze the system graphically. This equation is non-linear, but it is also separable

(6.6)

We can now work with just the left-hand side of this equation to obtain

(6.7)

With initial condition, $x(0) = x_0$, we can write down an equation for t in terms of x

(6.8)

but we want x in terms of t . This equation is hard to invert and is not particularly illuminating. We can make progress, however, by instead plotting \dot{x} versus x and analyzing some basic features of this phase portrait of the system.

Exercise 6.1

1. Sketch the phase portrait of this system, labeling stable and unstable fixed points.

2. Sketch $x(t)$ when $x_0 = \frac{\pi}{4}$.

6.5 Exponential Growth

We now consider some basic continuous growth models in population biology. Continuous growth models are especially suited for populations where reproduction can happen at any point in time.

For a discrete time model, we have:

$$(6.9)$$

meaning that the rate of increase is equal to the increase in the number of individuals $N(t + \Delta t) - N(t)$ divided by the initial number of individuals. What happens if we make the interval Δt shorter and shorter? For Δt going to zero, we obtain:

$$(6.10)$$

rearranging:

$$(6.11)$$

which is the differential equation for the exponential (Malthusian) growth you should be familiar with.

6.5.1 Analysis

Given the initial condition:

$$(6.12)$$

we can solve the equation. The equation is separable:

$$(6.13)$$

Integrating both sides:

$$(6.14)$$

$$(6.15)$$

$$(6.16)$$

$$(6.17)$$

$$(6.18)$$

$$(6.19)$$

The interpretation of the solution is the following: for $r > 0$ the population grows exponentially. For $r < 0$ the decay is exponential. For $r = 0$ the population is constant.

6.6 The Logistic Growth Model

6.6.1 Intraspecific Competition

We have considered a model in which populations grow without constraint. However, due to the finiteness of resources, an environment cannot sustain an infinite population. When resources start being scarce, a fierce competition among individuals of the same species begins. We call this effect intraspecific competition (i.e., within the same species).

To set the stage, let's examine some data from Gause (1932). To test the effects of intraspecific competition, Gause used brewer's yeast (*Saccharomyces cerevisiae*). His experiments were based on the growth of yeast in an environment with a limited nutrients. He started the population at some low level and then he measured the volume of the cells in time. The data is represented in Figure 6.1.

Clearly, yeast does not grow to infinity. Rather, after 1 day or so the population stops growing and remains around a given volume. Suppose that a population can grow up to a certain threshold K , usually known as carrying capacity. This model can be represented by the differential equation:

$$(6.20)$$

where r can be measured as the growth rate at very low density and K is the carrying capacity for the system, i.e., the density of individuals the habitat can support. Note that $r > 0$ and $K > 0$, otherwise it does not make biological sense (Why if $r < 0$ the equation does not make sense?)

6.6.2 Analysis

Plot $\frac{dN(t)}{dt}$ vs. $N(t)$ (see Figure 6.2). Using this graphical method we can count the number of equilibria and investigate their stability.

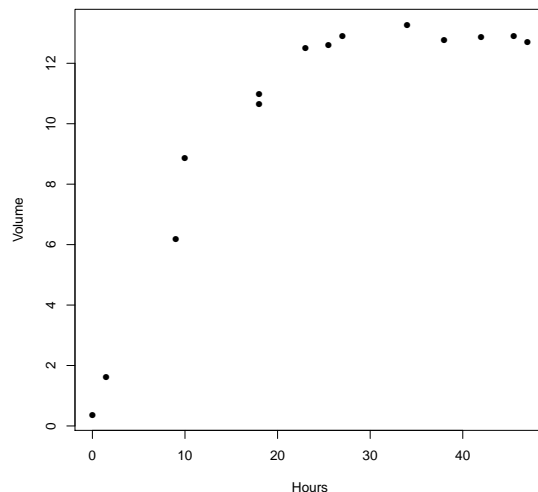


Figure 6.1 Data measured by Gause. Growth of brewer's yeast in an environment with limited nutrients.

We define an equilibrium point (or steady state, or fixed point or critical point) a point in which $\frac{dN}{dt} = 0$. In this point, the growth rate is zero: left alone, the population will not grow nor decay. In the logistic growth equation we have two fixed points: $N^* = 0$, $N^* = K$. This can be seen in the graph whenever $\frac{dN}{dt}$ crosses the zero line.

An equilibrium point N^* is stable if when we slightly perturb the system that is resting at N^* , $\lim_{t \rightarrow \infty} N(t) = N^*$. We can think about stability using the cartoon in Figure 6.3.

The red ball is at an unstable equilibrium point: slightly perturbing its position will set it in motion ending up at another point. The blue balls, on the other hand, will return to the same point once slightly perturbed: they are at stable equilibrium points.

For the logistic growth, we have that $N_1^* = 0$ is unstable, while $N_2^* = K$ is stable. If you start with zero bacteria, say, you will always have zero, but as soon as you have a single organism, you get exponential growth. Close to this fixed point, the system does not 'feel' the carrying-capacity. For single species dynamics, whenever growth rate is > 0 on the left of the fixed point and < 0 on the right of the fixed point we have stability. Therefore, if a species has density 0 and we slightly increase its density (decreasing it would be a biological nonsense) it will start growing until it reaches density K , when it will stop growing. It stops growing exponentially, and with the same rate that it exploded near 0.

We will now use some mathematics to get a qualitative understanding of equilibria. First, we want to find the equilibria. By definition, we want to solve the equation:

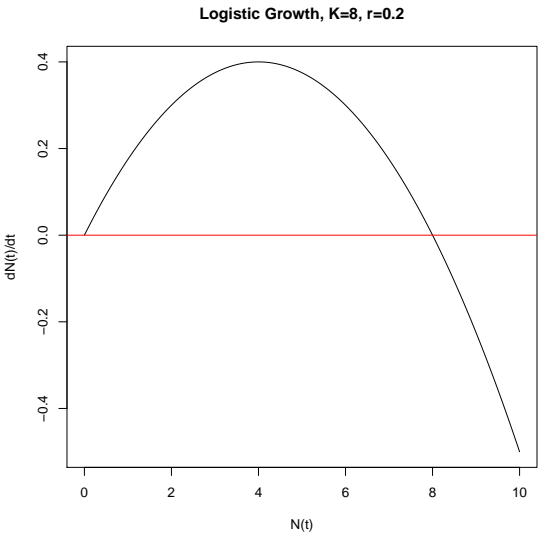


Figure 6.2 Graphical methods for logistic growth.

(6.21)

There is a trivial solution:

(6.22)

while the other solution can be found using:

(6.23)

(6.24)

(6.25)

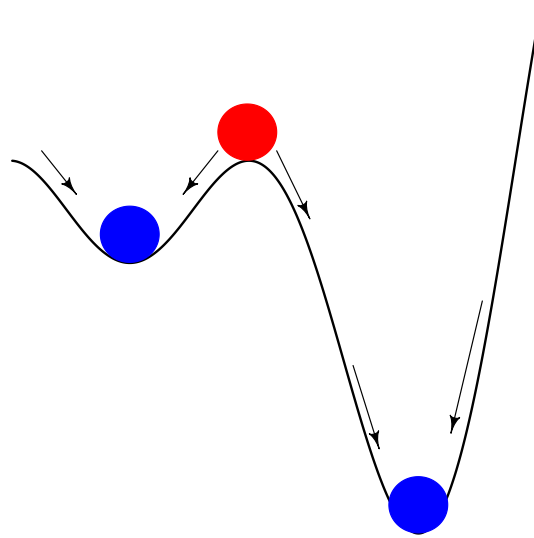


Figure 6.3 Stability as a metaphor.

We now want to measure the growth rate in the vicinity of an equilibrium. First, let us assume that $N(0) \ll K$ (i.e., the population at time 0 is at a much lower level than the carrying capacity). Then, $(1 - N(0)/K) \approx 1$. This means that the growth rate equation is approximately:

$$(6.26)$$

the population, when close to 0, will grow almost exponentially. In the same way, we can evaluate what happens when $N(0) \approx K$. In this case $(1 - N(0)/K) \approx 0$ and therefore:

$$(6.27)$$

the population will remain constant.

Finally, when $N(0) \gg K$ the term $(1 - N(0)/K)$ becomes negative (say the value is some unspecified $-y$). Then the approximate growth becomes:

$$(6.28)$$

the population decreases exponentially.

We want to know where the population grows at the fastest rate. To do so, we first take the derivative:

$$(6.29)$$

and equate it to zero:

$$(6.30)$$

$$(6.31)$$

we find that when the population is half-saturated, it grows at the fastest rate. Piecing these facts together, we can draw a qualitative solution for the differential equation (Figure 6.4).

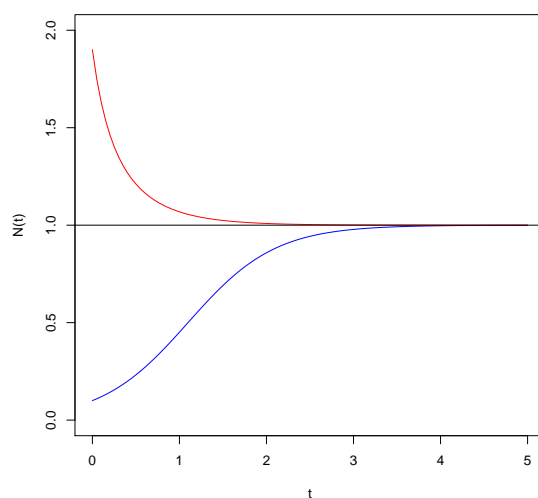


Figure 6.4 Logistic Growth: $K = 1$, $r = 2$. Blue: $N(0) = 0.1$, Red: $N(0) = 1.9$

This qualitative drawing captures the essential features of the solution:

(6.32)

6.7 Euler method for numerical integration

The Euler method is a very basic way to perform numerical integration on an ordinary differential equation. Suppose you wish to approximate a function, $\dot{x}(t) = f(x)$. You start at some initial condition, x_0 , and want to play the dynamics forward in time. The Euler method has a simple update rule, based on the definition of the derivative from your intro calculus course. For a small step in time, Δt , update x by

(6.33)

where x_n is an approximation to $x(t_n)$.

This method is very simple to code but suffers from the fact that you need to take step sizes, Δt , much smaller than any intrinsic time scale in your system to achieve good performance. Other numerical integration methods, such as Runge-Kutta, are included in most programming packages. You should simulate a system in which you know the solution to familiarize yourself with these numerical integration schemes.

6.8 Further reading

Would you like to know more about nonlinear dynamics and chaos? If so, a particularly accessible set of online lectures by Steve Strogatz is available on YouTube, and a link can be found in the **README** for this tutorial.

Tutorial 7

Introduction to UNIX

7.1 What is UNIX?

UNIX is an operating system (i.e., the software that let you interface with the computer) developed in the 1970s by a group of programmers at the AT&T Bell laboratories. Among them were Brian Kernighan and Dennis Ritchie, who also developed the programming language C. The new operating system was an immediate success in academic circles, with many scientists writing new programs to extend its features. This mix of commercial and academic interest led to the many variants of UNIX available today (e.g., OpenBSD, Sun Solaris, Apple OS X), collectively denoted as *nix systems. Most *nix systems are commercial, while Linux is the open source UNIX clone written from scratch by Linus Torvalds with the assistance of a loosely-knit team of hackers from across the internet. For this tutorial, we are going to focus on Linux and on Apple's OS X. Whenever needed, we will highlight the small differences between the two. Windows users can run **Cygwin**, which is a Linux emulator.

All *nix systems are multi-user and network-oriented, and store data as plain text files that can be exchanged between interconnected computer systems. Another characteristic is the use of a strictly hierarchical file system, discussed in Section 7.3.

7.2 Why use UNIX?

Many biologists are not familiar with coding in *nix systems, and – given that the learning curve is initially fairly steep – I start by listing the main advantages of these systems over possible alternatives.

First, UNIX is an operating system written by programmers for programmers. This means that UNIX is an ideal environment for developing your code and storing your data.

Second, in UNIX hundreds of small programs are available to perform simple tasks. These small programs can be stringed together efficiently, so that a single line of UNIX commands can perform complex operations—which otherwise would require writing a long

and complex program. The possibility of creating these pipelines for data analysis is especially important for biologists, as modern research groups produce large and complex data sets, whose analysis requires a level of automation that would be hard to achieve otherwise. For instance, imagine working with millions of files by having to open each one of them manually, or try opening your single 80Gb whole genome sequencing file in a software with a graphical user interface! In UNIX, you can simply string together a number of small programs, each performing a simple task, and create a complex pipeline that can be stored in a script (a text file containing all the commands). Then, you can let the computer analyze all of your data while you're having a cup of coffee.

Third, text is the rule: almost anything (including the screen, the mouse, etc.) in UNIX is stored in a text file, which means that all of your projects will be portable to other machines, and can be read and written without the need for sophisticated (and expensive) proprietary software, so that accessing your data does not depend on a specific software version or hardware requirement. Text files are (and always will be) supported by any operating system, so that you will be able to access your data decades from today (while this is not the case for most commercial software). The text-based nature of UNIX might seem unusual at first, especially if you are used to graphical interfaces and proprietary software. However, remember that UNIX has been around since the early 1970, and will likely be around at the end of your career. Thus, the hard work you are putting in learning UNIX will pay off over a lifetime.

The long history of UNIX means that a large body of tutorials and support web sites are readily available online.

Last but not least, UNIX is very stable, robust, secure, and—in the case of Linux—freely available.

7.3 Directory structure

In UNIX we speak of “directories”, while in a graphical environment the term “folder” is more common. These two terms are interchangeable, and refer to a structure that may contain sub-directories and files. The UNIX directory structure is organized hierarchically in a tree. As a biologist, you can think of this structure as a phylogenetic tree. The common ancestor of all directories is also called the “root” directory and is denoted by an individual slash (/). From the root directory, several important directories branch:

<code>/bin</code>	contains several basic programs.
<code>/etc</code>	contains configuration files.
<code>/dev</code>	contains the files connecting to devices such as the keyboard, mouse and screen.
<code>/home</code>	contains the home directory of each user (e.g., <code>/home/yourname</code> ; in OS X, your home directory is stored in <code>Users/yourname</code>).
<code>/tmp</code>	contains temporary files.

You will typically work in your home directory. From there, you can access the Desktop, Downloads, Documents, etc., directories you are likely familiar with. When you navigate the system, you are in one directory, and can move deeper in the tree, or upward towards the root. Section 7.5.2 discusses the commands that let you move between the hierarchical levels and determine your location within the directory structure.

7.4 Using the terminal

“Terminal” refers to the interface that you use to communicate with the kernel (the core of the operating system). The terminal is also called a “shell”, or command-line interface (CLI). It processes the commands you type, translates them for the kernel, and shows you the results of your operations. There are several shells available, and here we concentrate on the most popular one, the **bash** shell, which is the default shell in both Linux Ubuntu and OS X.

Ubuntu: To launch a shell in Ubuntu press **Ctrl + Alt + t**, or open the dash (hold the **Meta** key) and type “Terminal”. The shell will automatically start in your home directory `/home/yourname/`.

OSX: Open the Terminal.app, located in the folder “Applications → Utilities”. You can also search for it in spotlight by pressing the keys **Command + Space bar** and typing “Terminal”. The shell will open in you home directory `/Users/yourname/`.

For both Ubuntu and OS X, the command line prompt ends with a “dollar” (\$) sign. This means the terminal is ready to accept your commands. In these notes, a \$ sign at the beginning of a line of code signals that the command has to be executed in your terminal. You do not need to type the \$ sign in your terminal, just copy the command that follows it.

In UNIX, you can use the **Tab** key to reduce the amount you have to type, which in turn reduces errors caused by typos. When you press **Tab** in a (properly configured) shell, it will try to automatically complete your command, directory or file name (if multiple completions are possible, you can display them all by hitting the **Tab** key twice). Additionally, you can navigate the history of commands you typed by using the up/down arrows (you do not need re-type a command that you recently executed). There are also shortcuts that help pace through long lines of code:

```
Ctrl + A  Go to the beginning of the line
Ctrl + E  Go to the end of the line
Ctrl + L  Clear the screen
Ctrl + U  Clear the line before the cursor position
Ctrl + K  Clear the line after the cursor
Ctrl + C  Kill the command that is currently running
```

Ctrl + D Exit the current shell

Alt + F Move cursor forward one word (in OS X, **Esc + F**)

Alt + B Move cursor backward one word (in OS X, **Esc + B**)

Mastering these and other keyboard shortcuts will save you a lot of time. You may want to print this list and keep it next to your keyboard—in a while you will have them all memorized, and start using them automatically.

7.5 Basic UNIX commands

Here we introduce some of the most basic (and most useful) UNIX commands. We write the commands in **fixed-width font** and specific, user-provided input is capitalized in square brackets. Again, the brackets and special formatting are not required to execute a command in your terminal.

Many commands require some arguments (e.g., copy which file to where), and all can be modified using the several options available. Typically, options are either written as a dash followed by a single letter (older style, e.g., **-f**) or two dashes followed by words (newer style, e.g., **--full-name**).

7.5.1 How to get help in UNIX

UNIX ships with hundreds of commands. As such, it is impossible to remember them all, let alone remembering all the possible options. Fortunately, each command is described in detail in its manual page, which can be accessed directly from the shell by typing **man [COMMAND OF YOUR CHOICE]**. Use arrows to scroll up and down, and hit **q** to close the manual page. Checking the exact behavior of a command is especially important, given that the shell will execute any command you type without asking whether you know what you're doing (so that it will promptly remove all of your files, if that's the command you typed). You may be used to more forgiving (and slightly patronizing) operating systems, in which a pop-up window will warn you whenever something you're doing is considered dangerous. In UNIX, it is always better to consult the manual, rather than improvising.

7.5.2 Navigating the directory system

You can navigate the hierarchical UNIX directory system using these commands:

pwd print the path of the current working directory.

ls list the files and sub-directories in the current directory.

cd [NAMEOFDIR]
 change directory.

```
cd .. move one directory up
cd / move to the root directory
cd ~ move to your home directory.
```

7.5.3 Handling directories and files

Create and delete files or directories using the following commands:

```
cp [FROM] [TO]
    copy a file.
    The first argument is the file to copy. The second argument is where to copy
    it (either a directory or a file name).

mv [FROM] [TO]
    move or rename a file.
    Move a file by specifying two arguments: the file, and the destination directory.
    Rename a file by specifying the old and the new file name in the same directory.

touch [FILENAME]
    create an empty file.

rm [TOREMOVE]
    remove a file.
    rm -r allows to delete a directory recursively (i.e., including all files and sub-
    directories in it; use with caution!).

mkdir [DIRECTORY]
    make a directory.
    To create nested directories, use the option -p (e.g., mkdir -p d1/d2/d3).

rmdir [DIRECTORY]
    remove an empty directory.
```

7.5.4 Printing and modulating files

UNIX was especially designed to handle text files well, which is apparent when considering the multitude of commands dealing with text. Here are a few popular ones:

```
less [FILENAME]
    progressively print a file on the screen (press q to exit).1

cat [FILENAME]
    concatenate and print files.
```

¹Funny fact: there is a command called **more** that does the same thing, but with less flexibility. Clearly, in UNIX, **less** is **more**.

wc [FILENAME]
 word, line, character, and byte count of a file.

sort [FILENAME]
 sort the lines of a file and print the result to the screen.

uniq
 show only unique elements of an ordered list.

file [FILENAME]
 determine the type of a file.

head [FILENAME]
 print the head (i.e., first few lines of a file).

tail [FILENAME]
 print the tail (i.e., last few lines of a file).

diff [FILE1] [FILE2]
 show the differences between two files (can be installed in Cygwin).

7.5.5 Time and date

date
 print the current date.

cal
 display a calendar.

7.5.6 Miscellaneous

echo "[A STRING]"
 print the string [A STRING].

time
 time the execution of a command.

wget [URL]
 download the webpage at [URL] (available in Ubuntu, can be installed in OS X and Cygwin).

history
 lists the last commands you executed (10 by default).

7.6 Advanced UNIX commands

7.6.1 Redirection and pipes

So far, we have printed the output of each command (e.g., `ls`) directly to the screen. However, it is easy to direct the output to a file (“redirect”), or use it as the input of another command (“pipe”). Stringing commands together in pipes is the real power of UNIX—the ability to perform complex processing of large amounts of data in a single line of commands. First, we show how to redirect the output of a command into a file:

```
$ [COMMAND] > filename
```


Note that if the file `filename` exists, it will be overwritten. If instead we want to append to an existing file, we can use the `>>` symbol as in the following line:

```
$ [COMMAND] >> filename
```

When the command is very long and complex, we might want to redirect the content of a file as input to a command.

```
$ [COMMAND] < filename
```

To run a few examples, let's start by moving to our sandbox:

```
$ cd ~/QBio/Tutorials/T7-Intro_to_UNIX/Sandbox
```

(or `cd` to wherever you saved the `QBio` folder).

The command `echo` can be used to print a string on the screen. Instead of printing to the screen, we redirect the output to a file, effectively creating a file containing the string we want to print:

```
$ echo "My first line" > test.txt
```

We can see the result of our operation by printing the file to the screen using the command `cat`:

```
$ cat test.txt
```

To append a second line to the file, we use `>>`:

```
$ echo "My second line" >> test.txt
$ cat test.txt
```

We can redirect the output of any command to a file. For example, let's create a file listing all the files and directories accessible from the root of the directory system:

```
$ ls / >> ListRootDir.txt
$ cat ListRootDir.txt
```

Now, let's look at pipes (symbol `|`), which can be used to connect several commands.

Suppose we want to count how many files and directories are contained in the root directory. We can do this in several different ways, with or without pipes. A possible strategy would be to create a file containing all the names of the files and directories (as done above), and then use the command `wc -l` (count only the lines) to count them:

```
$ ls / > ListRootDir.txt
$ wc -l ListRootDir.txt
```

However, we can skip the creation of the file by simply piping the output of the command `ls` to the command `wc`:

```
$ ls / | wc -l
```

In the following sections, we are going to build increasingly long and complex pipelines. The idea is always to start with a command, and progressively add a piece after another to the pipeline, each time checking that the result is the desired one.

7.6.2 Selecting columns using cut

When dealing with tabular data, you will often encounter the Comma Separated Values (CSV) Standard File Format. The CSV format is platform and software independent, making it the standard output format of many experimental devices. The versatility of the file format should also make it your preferred choice when manually entering and storing data.

The main UNIX command you want to master for comma- or tab-delimited text files is `cut`. To show its main features, we will work with data on generation time of mammals published by Pacifici *et al.*, Nature Conservation, 2013. First, let's make sure we are in the right directory (QBio/Tutorials/T7-Intro.to_UNIX/Data), and then we can print the header (the first line, specifying the content of each column) of the CSV file:

```
$ head -n 1 Pacifici2013_data.csv
TaxID;Order;Family;Genus;Scientific_name;...
```

We now pipe the header to `cut`, specify the character to be used as delimiter, and extract the name of the first column, or the names of the first four columns:

```
$ head -n 1 Pacifici2013_data.csv | cut -d ';' -f 1
TaxID

$ head -n 1 Pacifici2013_data.csv | cut -d ';' -f 1-4
TaxID;Order;Family;Genus
```

In the next example, we work with the file content. We specify a delimiter, extract specific columns, and pipe the result to the `head` command—in order to display only the first few elements.

```
$ cut -d ';' -f 2 Pacifici2013_data.csv | head -n 5
Order
Rodentia
Rodentia
Rodentia
Macroscelidea

$ cut -d ';' -f 2,8 Pacifici2013_data.csv | head -n 3
Order;Max_longevity_d
Rodentia;292
```

```
Rodentia;456.25
```

In the next example, we specify the delimiter, extract the second column, skip the first line (the header) using the `tail -n +2` command (i.e., return the whole file starting from the second line), and finally display the first five entries:

```
$ cut -d ';' -f 2 Pacifici2013_data.csv | tail -n +2 | head -n 5
Rodentia
Rodentia
Rodentia
Macroscelidea
Rodentia
```

Now we pipe the result of the previous command to the `sort` command (which sorts the lines), and then again to `uniq`, which takes only the elements that are not repeated (entries need to be sorted for `uniq` to work properly). Effectively, we have created a pipeline to extract the names of all the Orders in the database, from Afrosoricida to Tubulidentata (a remarkable Order, which today contains only the aardvark).

```
$ cut -d ';' -f 2 Pacifici2013_data.csv | tail -n +2 | sort | uniq
Afrosoricida
Carnivora
Cetartiodactyl
...
```

This type of manipulation of delimited files is very fast and effective. It is an excellent idea to master the `cut` command, in order to start exploring large data sets without the need to open the file in a specialized programs.

7.6.3 Substituting characters using `tr`

We often want to substitute or remove a specific character in a text file (e.g., to convert a comma-separated file to a tab-separated file). Such a one-by-one substitution can be accomplished with the command `tr`. Let's look at some examples in which we use a pipe to pass a string to `tr`, which processes the text input according to the search term and specific options.

Substitute all characters `a` with `b`:

```
$ echo 'aaaabbbb' | tr 'a' 'b'
bbbbbbbb
```

Substitute every number in the range 1 through 5 with 0:

```
$ echo '123456789' | tr 1-5 0
000006789
```

Substitute **a** with 1, **c** with 2, and **d** with 3:

```
$ echo 'aabbccdde' | tr acd 123
11bb2233ee
```

We can also indicate ranges to substitute:

```
$ echo 'aabbccdde' | tr a-c 1-3
112233dde
```

Delete all occurrences of **a**:

```
$ echo 'aaaaabbbb' | tr -d a
bbbb
```

Squeeze all consecutive occurrences of **a**:

```
$ echo 'aaaaabbbb' | tr -s a
abbbb
```

Substitute spaces with tabs:

```
$ echo 'a b c d' | tr -s ' ' \t
a      b      c      d
```

Notice that the result is dependent on the options given to **tr**:

- s squeeze multiple, consecutive occurrences of the character listed in the search term.
- d delete all occurrences of the search term from the input.

Now we can apply the command **tr** and the commands we have showcased earlier, to create a new file, containing a subset of the data contained in **Pacifici2013.data.csv**, which we are going to use in the next Section.

First, we change directory to the sandbox:

```
$ cd ../Sandbox/
```

Because we were working in **Data**, moving one directory up (**..**) would bring us to **T7-Intro_to_UNIX**, from which we can move down to the **Sandbox**.

Now, we want to create a version of **Pacifici2013.data.csv**, containing only the **Order**, **Family**, **Genus**, **Scientific_name**, and **AdultBodyMass_g** (columns 2-6). Moreover, we want to remove the header, sort the lines according to body mass (with larger critters first), and have the values separated by spaces. This sounds like an awful lot of work, but we're going to see how this can be accomplished piping a few commands together.

First, let's remove the header:

```
$ cat ../Data/Pacifici2013_data.csv | tail -n +2
```

Then, take only the columns 2-6:

```
$ cat ../Data/Pacifici2013_data.csv | tail -n +2 | cut -d ';' -f 2-6
```

Now, substitute the current delimiter (;) by a space:

```
$ cat ../Data/Pacifici2013_data.csv | tail -n +2 | cut -d ';' -f 2-6  
| tr -s ';' ' '
```

To sort the lines according to body size, we need to exploit a few of the options for the command `sort`. First, we want to sort numbers, and thus we have to use the option `-n`. Second, we want larger values first, and thus need the option `-r`. Finally, we want to specify that the column we want to use to sort the data is the sixth, which can be accomplished using `-k 6`:

```
$ cat ../Data/Pacifici2013_data.csv | tail -n +2 | cut -d ';' -f 2-6  
| tr -s ';' ' ' | sort -r -n -k 6
```

That's it. We have created our first complex pipeline. To complete the task, we redirect the output of our pipeline to the file `BodyM.csv`.

```
$ cat ../Data/Pacifici2013_data.csv | tail -n +2 | cut -d ';' -f 2-6  
| tr -s ';' ' ' | sort -r -n -k 6 > BodyM.csv
```

One might object that the same task could have been accomplished with a few clicks by opening the file in a spreadsheet editor. However, suppose you have to repeat this task many times, for instance to reformat every file that is produced by a laboratory device. Then it is convenient to automate this task, such that it can be run with a single command. This is exactly what we are going to do in Section 7.7.

7.6.4 Selecting lines using `grep`

`grep` is a powerful command that finds all the lines of a file that match a given pattern. You can return or count all occurrences of the pattern in a large text file without ever opening it. `grep` is based on the concept of regular expressions.

We will test the basic features of `grep` using the file we just created in section 7.6.3. The file contains data on thousands of species:

```
$ wc -l BodyM.csv  
5426 BodyM.csv
```

Let's see how many wombats (family `Vombatidae`) are contained in the data. First we display the lines that contain the term "Vombatidae":

```
$ grep Vombatidae BodyM.csv  
Diprotodontia Vombatidae Lasiorhinus Lasiorhinus krefftii 31849.99  
Diprotodontia Vombatidae Lasiorhinus Lasiorhinus latifrons 26163.8  
Diprotodontia Vombatidae Vombatus Vombatus ursinus 26000
```

Now we add the option `-c` to count:

```
$ grep -c Vombatidae BodyM.csv
3
```

Next, we have a look at the genus *Bos* in the data file:

```
$ grep Bos BodyM.csv
Cetartiodactyla Bovidae Bos Bos sauveli 791321.8
Cetartiodactyla Bovidae Bos Bos gaurus 721000
Cetartiodactyla Bovidae Bos Bos mutus 650000
Cetartiodactyla Bovidae Bos Bos javanicus 635974.3
Cetartiodactyla Bovidae Boselaphus Boselaphus tragocamelus 182253
```

Besides all the members of the *Bos* genus, we also match one member of the genus *Boselaphus*. To exclude it, we can use the option `-w`, which prompts `grep` to match only full words:

```
$ grep -w Bos BodyM.csv
Cetartiodactyla Bovidae Bos Bos sauveli 791321.8
Cetartiodactyla Bovidae Bos Bos gaurus 721000
Cetartiodactyla Bovidae Bos Bos mutus 650000
Cetartiodactyla Bovidae Bos Bos javanicus 635974.3
```

Using the option `-i` we can search case-insensitive:

```
$ grep -i Bos BodyM.csv
Proboscidea Elephantidae Loxodonta Loxodonta africana 3824540
Proboscidea Elephantidae Elephas Elephas maximus 3269794
Cetartiodactyla Bovidae Bos Bos sauveli 791321.8
Cetartiodactyla Bovidae Bos Bos gaurus 721000
...
```

Sometimes, we want to know which lines precede or follow the one we want to match. For example, suppose we want to know which mammals have body weight most similar to the gorilla (*Gorilla gorilla*). The species are already ordered by size (see 7.6.3) thus we can now simply print the two lines before the match using the option `-B 2`, and the two lines after the match using `-A 2`:

```
$ grep -A 2 -B 2 "Gorilla gorilla" BodyM.csv
Cetartiodactyla Bovidae Ovis Ovis ammon 113998.7
Cetartiodactyla Delphinidae Lissodelphis Lissodelphis borealis
113000
Primates Hominidae Gorilla Gorilla gorilla 112589
Cetartiodactyla Cervidae Blastocerus Blastocerus dichotomus 112518.5
Cetartiodactyla Iniidae Lipotes Lipotes vexillifer 112138.3
```

Use option `-n` to show the line number of the match. For example, the gorilla is the 164th largest mammal in the database:

```
$ grep -n "Gorilla gorilla" BodyM.csv
164:Primates Hominidae Gorilla Gorilla gorilla 112589
```

To print all the lines that do not match a given pattern, use the option `-v`. For example, to get the other species of the genus *Gorilla* with the exception of *Gorilla gorilla*, we can use:

```
$ grep -n Gorilla BodyM.csv | grep -v gorilla
143:Primates Hominidae Gorilla Gorilla beringei 149325.2
```

To match one of several strings, use `grep "string1\|string2"`:

```
$ grep -w "Gorilla\|Pan" BodyM.csv
Primates Hominidae Gorilla Gorilla beringei 149325.2
Primates Hominidae Gorilla Gorilla gorilla 112589
Primates Hominidae Pan Pan troglodytes 45000
Primates Hominidae Pan Pan paniscus 35119.95
```

You can use `grep` on multiple files at a time! Simply, list all the files to use instead of just one file. Finally, use the recursive search option `-r` to search for patterns within all the files in a directory, for example:

```
$ grep -r "Gorilla" ../Data/
```

7.7 Basic scripting

Once a pipeline is in place, it is easy to make it into a “script”, a collection of commands that goes through the pipeline in an automated manner. To show how this can be accomplished, we are going to turn the pipeline in Section 7.6.3 into a script.

First, we need to create a file for our UNIX script, which we can edit using a text editor. The typical extension for such a shell file is `.sh`. In this example we want to create the file `ExtractBodyM.sh`, which we can open using our favorite text editor.

Create the empty file:

```
$ touch ExtractBodyM.sh
```

and open it with a text editor (e.g., `gedit` in Ubuntu or `emacs` in OS X):

```
$ gedit ExtractBodyM.sh &
```

the “ampersand” (`&`) at the end of the line prompts the terminal to open the editor in the background, so that you can still use the same shell while working on the file.

Now copy the pipeline into the file:

```
cat ../Data/Pacifici2013_data.csv | tail -n +2 | cut -d ';' -f 2-6 | tr
-s ';' ' ' | sort -r -n -k 6 > BodyM.csv
```

and save the file. Now you can run the script, by calling the command **bash**:

```
$ bash ExtractBodyM.sh
```

It is a great idea to immediately comment the script to help you remember what the script does when you return to it later. You can add comments using the “hashtag” sign (#):

```
1 # Take a csv file delimited by ';'
2 # Remove the header
  # Make space separated
  # Sort according to the 6th (numeric) column in descending order
5 # Redirect to a file
  cat ../Data/Pacifici2013_data.csv | tail -n +2 | cut -d ';' -f 2-6 | tr
    -s ';' ' ' | sort -r -n -k 6 > BodyM.csv
```

This script is very specific: it only converts the file `../Data/Pacifici2013_data.csv` into `BodyM.csv`, while it would be better to leave these names to be decided by the user, so that the script can be called on any file with the same format. This is very easy to accomplish within the **bash** shell, as one can use generic arguments (i.e., variables), indicated by the dollar sign, followed by the argument number:

```
1 # Take a csv file delimited by ';' (First argument)
2 # Remove the header
  # Make space separated
  # Sort according to the 6th (numeric) column in descending order
5 # Redirect to a file (Second argument)
  cat $1 | tail -n +2 | cut -d ';' -f 2-6 | tr -s ';' ' ' | sort -r -n -k
    6 > $2
```

Now you can launch the modified script by specifying the input and output files from the command line:

```
$ bash ExtractBodyM.sh ../Data/Pacifici2013_data.csv BodyM.csv
```

The final step is to make the script directly executable, so that you can skip the preceding **bash**. We can do so by changing the permissions of the file:

```
$ chmod +x ExtractBodyM.sh
```

and adding a special line at the beginning of the file, telling UNIX where to find the program (here **bash**) to execute the script:


```

#!/bin/bash
2
# The previous line is not a comment, but rather a special line
# telling UNIX where to find the program to execute the string.
5 # It should be your first line in all bash scripts

# What the script does:
8 # Take a csv file delimited by ';' (First argument)
# Remove the header
# Make space separated
11 # Sort according to the 6th (numeric) column in descending order
# Redirect to a file (Second argument)
cat $1 | tail -n +2 | cut -d ';' -f 2-6 | tr -s ';' ' ' | sort -r -n -k
    6 > $2

```

which can be invoked as:

```
$ ./ExtractBodyM.sh ../Data/Pacifici2013_data.csv BodyM.csv
```

Note the `./` in front of the script's name in order to execute the file.

Such a long UNIX pipe can be complicated to read and understand at a later point, in which case it is better to break it into pieces and save the individual output of each part into a temporary file that can be cleaned up:

```

#!/bin/bash
2 # What the script does:
# Take a csv file delimited by ';' (First argument)
# Remove the header
5 # Make space separated
# Sort according to the 6th (numeric) column in descending order
# Redirect to a file (Second argument)
8
# remove the header
cat $1 | tail -n +2 > $1.tmp1
11 # extract columns
cat $1.tmp1 | cut -d ';' -f 2-6 > $1.tmp2
# make space-separated
14 cat $1.tmp2 | tr -s ';' ' ' > $1.tmp3
# sort and redirect to output
cat $1.tmp3 | sort -r -n -k 6 > $2
17 # remove temporary, intermediate files
rm $1.tmp*

```

which is much more readable (although a little more wasteful, as it creates temporary files only to then delete them), and easier to “debug” (just comment the last line out and inspect the temporary files).

Tutorial 8

Data Visualization

8.1 Overview

In this tutorial, we will discuss some basic guidelines for displaying data. Next, you will apply these concepts to make improvements on a figure of your own or create a new figure with the data provided. Most of the material for this tutorial will be presented in slides that you can download from the course repository later.

8.2 Basic guidelines for displaying data

The first step in making an effective display of your data is to decide how to best summarize, collapse, or slice through your raw data to convey the quantitative result you are discussing in the text of your paper or proposal. In some sense, displaying your data properly starts with good writing. As a good first read on this topic, you will find an article by Gopen and Swan entitled, ‘The science of scientific writing,’ in the Readings section of this tutorial. Once you have decided what data you would like to display, here are a few simple rules to follow to make your figure clear and effective:

- Eliminate any uninformative or unnecessary distractions from your plot. This may include background texture or shading and data panels that do not communicate the points you are making in the text.
- Avoid gratuitous color in your plots.
- Make sure your axes labels and all other text in your plots is readable. A general rule of thumb is that text in publication figures should be between 6 and 12 pt.
- Avoid placing text over low contrast or textured backgrounds.
- Whenever possible, avoid pie charts.

- Whenever possible, avoid spaghetti plots.

8.3 An example dataset and published plot

To compare your own best figure with a recent Economist graphic, you may download the UN world population projection data from the Data section of this tutorial: use `UN population forecasts from Economist charts - Country.csv` or `UN population forecasts from Economist charts - Region.csv`. The Economist's take on these data is shown on the following page.

Exercise 8.1

1. Open up the data in these files and explore it. What to you is the most compelling feature in these data?
2. Do you think the Economist graphic captures that aspect of the data?
3. Make a plot to highlight the feature that most intrigued you.
4. Write a concise caption to explain the data shown in your plot.

8.4 A brief introduction to ggplot

You were first introduced to the 'grammar of graphics' in the programming tutorials. We will go over some basic commands here, as well. You will find some useful short guides to `ggplot2` in the Readings section for this tutorial.

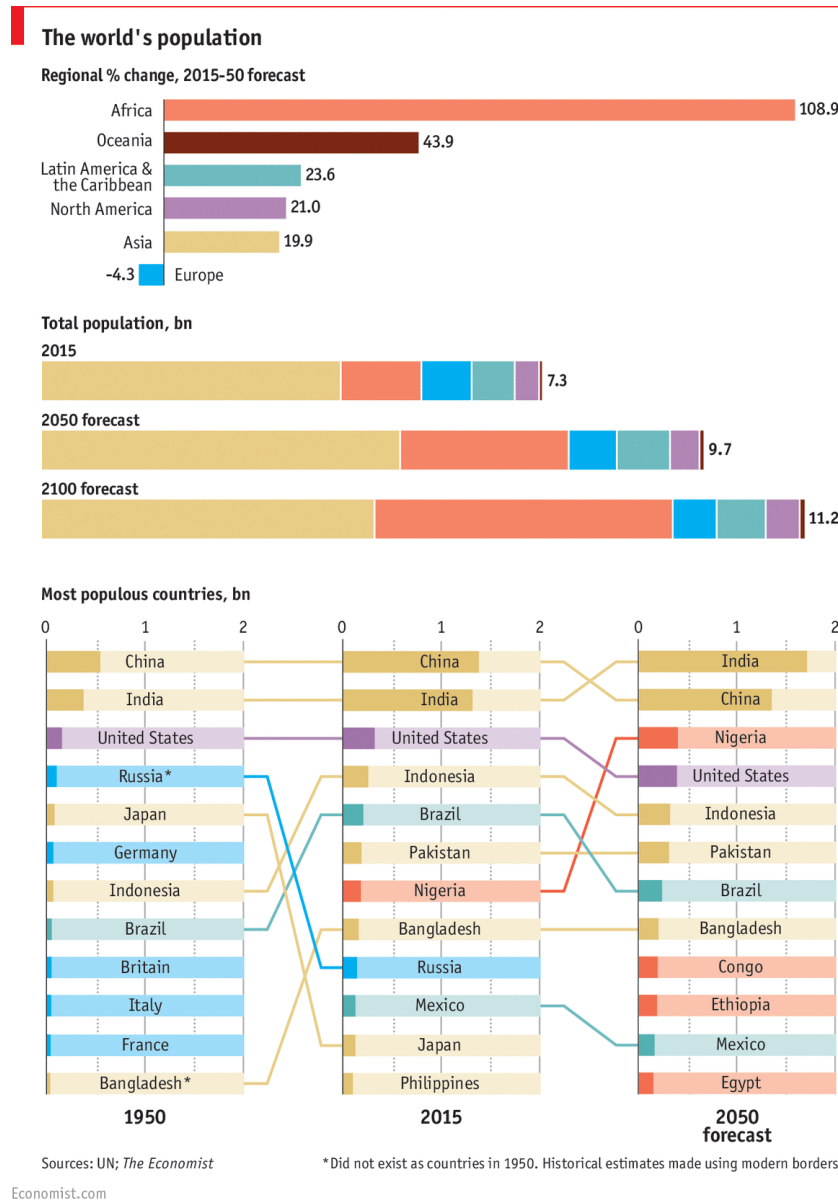


Figure 8.1 A graphic from the August 4, 2015 issue of the Economist, displaying global population forecasts using UN data tables. Highlights of the large data set are displayed here. A link to the original figure and caption is included in the README for this tutorial. <http://www.economist.com/blogs/graphicdetail/2015/08/daily-chart-growth-areas>