

R: An Introductory Session

The following session is intended to introduce to you some features of the R environment: the best way to learn R is by using it. Many features of the system will be unfamiliar and puzzling at first, but this will soon disappear. Also, although you may not yet understand all the statistical details, the plan is to type the commands and see what happens as a result.

The left column gives commands; the right column gives brief explanations and suggestions.

Login to R by clicking on the `intro.RData` file after you have downloaded it from MyUni.

<code>ls()</code>	Obtain a listing of all objects currently available to your session.
<code>library(MASS)</code>	This command makes the datasets from Venables and Ripley available.

<code>help.start()</code>	Start the on-line help facility. Alternatively, use the help drop down window as described in the Computing Guide.
<code>help(help)</code> <code>?help</code>	This is another way to get help. Read the help page about how to use help.

<code>cystfibr</code>	Typing the name of an object causes it to be displayed. (This is not a good idea for very large datasets.) The <code>cystfibr</code> data frame has 25 rows and 10 columns; it contains lung function data for cystic fibrosis patients 7–23 years old.
<code>is.data.frame(cystfibr)</code> <code>length(cystfibr)</code> <code>attributes(cystfibr)</code>	We determine the type of an object and its attributes (names of variables or components, etc) by using commands like this.
<code>cystfibr\$fev1</code>	The \$ symbol allows us to refer to individual variables without ‘attaching’ the dataframe. <code>fev1</code> is a numeric vector of measurements on forced expiratory volume.

Some graphical features in R:

<code>x <- rnorm(1000)</code> <code>y <- rnorm(1000)</code>	Generate 1000 pairs of normal variates.
<code>truehist(c(x, y+3), nbins=25)</code>	Histogram of a mixture of normal distributions. Experiment with the number of bins (25) and the shift (3) of the second component.
<code>?truehist</code>	Read about the optional arguments.
<code>contour(dd <- kde2d(x,y))</code>	A two-dimensional density plot.

<code>image(dd)</code>	A pseudo-colour plot.
<code>hull <- chull(x, y)</code>	Find the convex hull of x and y in the plane.
<code>plot(x, y)</code>	Plot the points in the plane, and mark in their convex hull.
<code>polygon(x[hull], y[hull], dens=15)</code>	
<hr/>	
<code>data(hills)</code>	This dataset is from the MASS library and contains data on record times of Scottish hill races against distance and total height climbed.
<code>hills</code>	List the data.
<code>pairs(hills)</code>	Show a matrix of pairwise scatterplots.
<code>objects()</code>	See which R objects you have created.
<code>rm(x, y, hull)</code>	Remove objects no longer needed (i.e., cleanup).
<hr/>	

Some basic data manipulation in R:

<code>> x <- 1</code>	Arithmetic expressions typed into R are evaluated and printed out. Try these examples (the R prompt is included here for clarity).
<code>> y <- 2</code>	
<code>> x+y</code>	
<code>> log(x/(1+3*y))</code>	
<code>> w <- x+y</code>	
<code>> z <- log(x/(1+3*y))</code>	To save the result of a calculation as a new R object, use the assignment operator <code><-</code> . Try these examples.
<code>> w</code>	
<code>> z</code>	
<hr/>	

<code>small.data</code>	This object is a list containing four components, each of which is a matrix. This structure is typical of microarray data and these are data from 3 slides (each slide represents a different experiment) for 10 genes. Examine the attributes etc of <code>small.data</code> . The data objects in an arithmetic expression can be a combination of scalars, vectors or matrices. If they are all matrices of the same dimension, then the expression is evaluated for each component.
<code>small.data\$R</code>	These are matrices of red foreground and background intensities.
<code>small.data\$Rb</code>	
<code>Red <- small.data\$R - small.data\$Rb</code>	We subtract the backgrounds to obtain the background-adjusted intensities.
<code>Green <- small.data\$G - small.data\$Gb</code>	Obtain the green background-adjusted intensities and find the log base 2 intensity ratio M .
<code>M <- log(Red/Green, base=2)</code>	
<code>M</code>	

As an exercise, calculate the matrix A of average log intensities

$$A = (\log_2 \text{Red} + \log_2 \text{Green})/2.$$

We will next look more closely at the cystic fibrosis data.

<pre>attach(cystfibr)</pre>	This makes the columns (i.e., variables) available by name.
<pre>par(mex=0.5)</pre>	Obtain pairwise scatterplots between all the variables in the dataset. The arguments <code>gap</code> and <code>cex.labels</code> control the visual appearance by removing the space between subplots and decreasing the font size; <code>mex</code> reduces the interline distance in the margins. <code>plot(cystfibr)</code> gives a similar plot.
<pre>pairs(cystfibr, gap=0, cex.labels=0.9)</pre>	
	<i>Which variables appear to be closely related?</i>
<pre>fm <- lm(pemax~age) summary(fm)</pre>	Fit a simple linear regression of <code>pemax</code> on <code>age</code> (in years) and look at the analysis; <code>pemax</code> is a numeric vector of maximum expiratory pressure.
<pre>plot(age, pemax, col="blue")</pre>	Make a standard scatterplot and colour the points blue.
<pre>title("Max expiratory pressure vs age")</pre>	Add a title to the plot.
<pre>abline(fm, lty=3, col="red")</pre>	Add in the red regression line using a different line type; <code>abline()</code> is able to extract the information it needs from the fitted regression object.
<pre>lines(lowess(age,pemax), col="green")</pre>	Fit a smooth regression curve using a modern regression function, and plot the smoothed line. Use the help facility to learn more about <code>lowess</code> .
<pre>plot(fitted(fm), resid(fm), xlab="Fitted values", ylab="Residuals", main= "Residuals vs Fitted")</pre>	A standard regression diagnostic plot to check for heteroscedasticity (i.e., unequal variance). Is there any evidence of heteroscedasticity here? Guess what all the arguments do.
<pre>qqnorm(resid(fm), main= "Q-Q Plot")</pre>	A normal scores plot to check for skewness, kurtosis and outliers.
<pre>qqline(resid(fm), col="green")</pre>	At any time you can make a hardcopy of the graphics window by clicking on the appropriate box of the window and selecting the <code>Print</code> option.
<pre>detach(cystfibr)</pre>	Remove the <code>cystfibr</code> data frame from the search list.
<pre>rm(fm)</pre>	Clean up again.

Next we will look at using R for weighted regression with an exercise using simulated data.

<code>x <- seq(1, 20, 0.5)</code>	Make $x = (1, 1.5, \dots, 19.5, 20)$ and list it.
<code>x</code>	
<code>w <- 1 + x/2</code>	<code>w</code> will be used as a 'weight' vector and to give the standard deviations of the errors.
<code>y <- x + w*rnorm(x)</code>	
<code>dum <- data.frame(x,y,w)</code>	Make a data frame of three columns named <code>x</code> , <code>y</code> and <code>w</code> , and look at it. Remove the original <code>x</code> , <code>y</code> and <code>w</code> .
<code>dum</code>	
<code>rm(x,y,w)</code>	
<code>fm <- lm(y~x, data=dum)</code>	Fit a simple linear regression of <code>y</code> on <code>x</code> and look at the analysis.
<code>summary(fm)</code>	
<code>fm1 <- lm(y~x, data=dum, weight=1/w^2)</code>	Since we know the standard deviations, we can do a weighted regression.
<code>summary(fm1)</code>	
<code>library(modreg)</code>	This is a library of nonparametric smoothing functions, including splines.
<code>lrf <- loess(y~x, dum)</code>	Fit a smooth (nonparametric) regression curve. This is like fitting many regressions over the x - range and 'linking' them together. <code>loess</code> is a newer function than <code>lowess</code> .
<code>attach(dum)</code>	
<code>plot(x,y)</code>	
<code>lines(spline(x, fitted(lrf)), col=2)</code>	Add in the local regression curve using a spline interpolation between the calculated points.
<code>abline(0,1,lty=3,col=3)</code>	Add in the true regression line (intercept 0, slope 1) with a different line type and colour.
<code>abline(fm, col=4)</code>	
<code>abline(fm1, lty=4,col=5)</code>	Finally add in the weighted regression line. This one should be the most accurate estimate, but may not be of course.
<code>detach(dum)</code>	Clean up again.
	An an exercise, obtain the standard diagnostic plots as before.
<code>q()</code>	Quit R.

Patty Solomon

July 2004