

Learning R

Gavin L. Simpson

February, 2017

Introduction to R

Why R?

- R was designed from the ground up as a language for data analysis
- It is free and open source, and available on all the major OSes
- Huge package ecosystem (> 10,000) covering bewildering array of statistical methods, data visualizations, data import and manipulation
- Cutting edge; R is used by thousands of statisticians & R code or a package often accompanies papers developing new methods
- For the most part, a great community providing help, blog posts etc

What is R?

- The S statistical language was started at Bell Labs on May 5, 1976
- A system for general data analysis jobs that could replace the *ad hoc* creation of FORTRAN applications
- The S language was licensed by Insightful Corporation for use in their *S-PLUS* software
- In 2004 Insightful bought the S language from Lucent (formerly AT&T and before that Bell Labs)
- Robert Gentleman and Ross Ihaka designed a language that was compatible with S but which worked in a different way internally
- They called this language R
- There was a lot of interest in R and eventually it was made Open Source under the GNU GPL-2
- R has drawn around it a group of dedicated stewards of the R software — **R Core**

- The R homepage is located at: www.r-project.org
- The download site is called CRAN — the **Comprehensive R Archive Network**
- CRAN is a series of mirrored web servers to spread the load of thousands of users downloading R and associated packages
- The CRAN master is at: cran.r-project.org

Starting R and other preliminaries

- You start R in a variety of ways depending on your OS
- R starts in a **working directory** where it looks for files and saves objects
- Best to run R in a new directory for each project or analysis task
- `getwd()` and `setwd()` get and set the working directory
- To exit R, the function `q()` is used
- You will be asked if you want to save your workspace; invariably you should answer `n` to this

- R comes with a lot of documentation
- To get help on functions or concepts within R, use the "?" operator
- For help on the `getwd()` function use: `?getwd`
- Function `help.search("foo")` will search through all packages installed for help pages with "foo" in them
- How the help is displayed is system dependent
- Google is your friend
- StackOverflow's R tag: stackoverflow.com/questions/tagged/r

- Type commands at prompt ">" and these are evaluated when you hit RETURN
- If a line is not syntactically complete, the prompt is changed to "+ "
- If returned object not assigned, it is printed to console
- Assigning the results of a function call achieved by the assignment operator "<-"
- Whatever is on the right of "<-" is assigned to the object named on the left of "<-"
- Enter the name of an object and hit RETURN to print the contents
- `ls()` returns a list of objects currently in your workspace

Working with R & entering commands

```
> 5 * 3
```

```
[1] 15
```

```
> radius <- 5
```

```
> pi * radius^2
```

```
[1] 78.53982
```

```
> ans <- 5 * 3
```

```
> ans
```

```
[1] 15
```

```
> ans2 <- ans + 20
```

```
> ans2
```

```
[1] 35
```

```
> ls()
```

```
[1] "ans"    "ans2"   "radius"
```

Data structures

The basic data structures

Following Hadley Wickham's *Advanced R*, the basic data structures in R can be classified according to

1. whether their contents are all the same (homogeneous) or not,
2. the dimensionality of the object

Homo	geneous Hete	rogeneous
1d	Atomic vector	List
2d	Matrix	Data frame
nd	Array	NA

The best way to understand which data structure an R object is, is to use `str()`

Vectors are the basic type of data object in R and there are two main types

1. *atomic* vectors
2. lists

Each with three common properties

- What type of vector the object is: `typeof()`
- Its length: `length()`
- Attributes, which are extra information or metadata: `attributes()`, `attr()`

Atomic vectors differ fundamentally from lists because atomic vectors can only contain elements that are *all of the same type*

The four main types of atomic vector are

1. logical
2. integer
3. double
4. character

Two other types are less often encountered; raw and complex vectors

We create atomic vector with `c()`, the *concatenate* or *combine* function

```
> dbl <- c(1, 2, 3)
> int <- c(1L, 2L, 3L)
> logi <- c(TRUE, FALSE, TRUE)      # Avoid using c(T, F, T)
> chr <- c("Hello", "World")
```

If an element (observation) is missing, use **NA**

All elements of atomic vectors must be of the same type. If you mix types or attempt to combine atomic vectors of different types, they will be *coerced* towards the most general type.

The ordering is (from least to most general)

1. logical
2. integer
3. double
4. character

```
> str(c("a", 1))  
chr [1:2] "a" "1"  
  
> str(c(TRUE, 10))  
num [1:2] 1 10  
  
> str(c(10L, 2))  
num [1:2] 10 2
```

One handy coercion is the conversion of logical vectors to numeric (integer or double).

A **TRUE** is **1**, whilst a **FALSE** is **0**, which allows them to be used in numeric mathematical operations

Coercion happens atomically, but you can control this by explicitly coercing to the required type with

- `as.character()`
- `as.double()`
- `as.integer()`
- `as.logical()`

```
> x <- c(FALSE, FALSE, TRUE)
> as.numeric(x)
```

```
[1] 0 0 1
```

```
> sum(x)                                # count up the TRUES
```

```
[1] 1
```

```
> mean(x)                               # proportion of TRUE
```

```
[1] 0.3333333
```


Lists are exceedingly common in R — it's often how fitted statistical models are stored

Lists are *general* vectors because the elements they contain can be of any type — lists can even contain lists

```
> x <- list(1:3, "a", c(TRUE, FALSE, TRUE), c(2.3, 5.9))
> str(x)
```

```
List of 4
 $ : int [1:3] 1 2 3
 $ : chr "a"
 $ : logi [1:3] TRUE FALSE TRUE
 $ : num [1:2] 2.3 5.9
```

```
> x <- list(list(list(list)))
> str(x)
```

```
List of 1
 $ :List of 1
  ..$ :List of 1
    .. ..$ :function (...)
```

Factors

Factors are useful for storing data where observations can take on one of a finite set of values

Factors combine integer vectors with attributes to create a new data structure

They have `class()` "factor" and a `levels()` attribute which lists the set of values the elements can take

```
> x <- factor(c("x","y","z","z","x","z","y"))
> x

[1] x y z z x z y
Levels: x y z

> class(x)

[1] "factor"

> levels(x)

[1] "x" "y" "z"

> x[2] <- "c"                                     # throws a warning

Warning in `[<-factor`(`*tmp*`, 2, value = "c"): invalid factor level, NA
generated

> x

[1] x <NA> z z x z y
Levels: x y z

> x <- factor(c("x","y","z","z"), levels = c("z","y","x")) # specify levels
```

What separates atomic vectors from arrays is the presence of a `dim()` attribute

As arrays are really atomic vectors with this extra attribute, they can only contain elements of a single type

Matrices are a special case of a multi-dimensional array, where there are only two dimensions

```
> m <- matrix(1:6, ncol = 3, nrow = 2)
```

```
> m
```

```
      [,1] [,2] [,3]  
[1,]    1    3    5  
[2,]    2    4    6
```

```
> m <- matrix(1:6, c(3,2))
```

```
>
```

```
> m <- 1:6
```

```
> dim(m) <- c(3,2)
```

```
> dim(m) <- c(2,3)
```

```
> m
```

```
      [,1] [,2] [,3]  
[1,]    1    3    5  
[2,]    2    4    6
```

Data frames

Data frames are a bit like Excel worksheets for R; they are like 2-d matrices with the exception that the columns can store different types of objects

Internally, data frames are lists, with the extra restriction that each component (column) of the data frame has to be of the same length

Data frame can be created using `data.frame()`

```
> df <- data.frame(x = 1:3, y = c("a", "b", "c"))
> str(df)

'data.frame':  3 obs. of  2 variables:
 $ x: int  1 2 3
 $ y: Factor w/ 3 levels "a","b","c": 1 2 3

> df <- data.frame(x = 1:3, y = c("a", "b", "c"), stringsAsFactors = FALSE)
> is.data.frame(df)

[1] TRUE

> nrow(df)

[1] 3

> class(df)

[1] "data.frame"
```

Additional columns and rows can be added to a data frame using `cbind()` and `rbind()` respectively

```
> cbind(df, data.frame(z = 3:1))
```

```
  x y z  
1 1 a 3  
2 2 b 2  
3 3 c 1
```

```
> rbind(df, data.frame(x = 10, y = "z"))
```

```
  x y  
1 1 a  
2 2 b  
3 3 c  
4 10 z
```

Subsetting is an incredibly useful and powerful way of manipulating data objects — but it is hard to learn initially as there is a lot of detail to remember

Subsetting atomic vectors

The main subsetting function is `[]` — it takes one of several types of arguments which determines how the vector is subset

- *positive integers* return the specified elements
- *negative integers* omit the specified elements from the result
- *logical vectors* return elements where the logical vector is **TRUE**
- *nothing* returns the original vector
- *zero* returns a zero-length vector

Subsetting a list works the same way with `[]`

```
> x <- c(1.3, 4.5, 2.3, 4.2, 5.4)
> x[c(3, 1)]

[1] 2.3 1.3

> x[-c(3, 1)]

[1] 4.5 4.2 5.4

> x[x > 3]

[1] 4.5 4.2 5.4

> x[c(TRUE, FALSE)]

[1] 1.3 2.3 5.4

> x[]

[1] 1.3 4.5 2.3 4.2 5.4

> x[0]

numeric(0)

> (y <- setNames(x, letters[1:5]))

  a    b    c    d    e 
1.3 4.5 2.3 4.2 5.4 

> y[c("a", "c")]

  a    c 
1.3 2.3
```

The main subsetting function is `[` — it takes one of several types of arguments which determines how the array is subset

- a pair of 1-d vectors, one for each dimension
- a single 1-d vector
- a matrix

```
> a <- matrix(1:9, nrow = 3)
> colnames(a) <- LETTERS[1:3]
> a[1:2, ]

      A B C
[1,] 1 4 7
[2,] 2 5 8
> a[, 2]

[1] 4 5 6
```


Subsetting data frames

Again, the main subsetting function is `[` — you can use it to subset a data frame as if it were

- a matrix, with an indexing vector for the rows and the columns
- a list

```
> df <- data.frame(x = 1:3, y = 3:1, z = letters[1:3])
> df[df$x == 2, ]

  x y z
2 2 2 b

> df[c("x", "y")]

  x y
1 1 3
2 2 2
3 3 1

> df[, c("x", "y")]

  x y
1 1 3
2 2 2
3 3 1

> str(df["x"])

'data.frame':  3 obs. of  1 variable:
 $ x: int  1 2 3

> str(df[, "x"])                                # drops the empty dimension

int [1:3] 1 2 3
```

Plotting

R is good at plotting; in-built functionality for producing publication-ready figures in a range of formats

Three main plotting paradigms in the R ecosystem

1. Base graphics
2. **lattice** package
3. **ggplot2** package

The latter, **lattice** and **ggplot**, both use the underlying **grid** primitive graphics toolkit

In general, R's graphics are like drawing with pen on paper; once you draw anything that sheet of paper is no-longer pristine and you can't erase anything you have drawn

Plotting using base graphics

Standard plotting command is `plot()`

Takes one or two arguments of coordinates

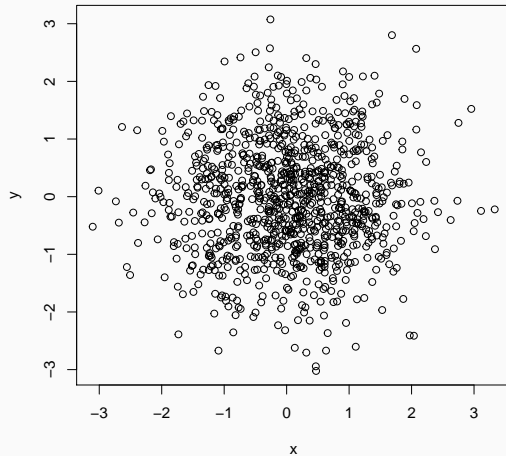
By default draws a scatterplot

Vast array of parameters to alter look of plots; see

`?par`

Best results often come from building a plot up from bits

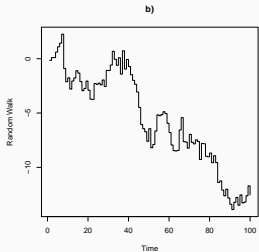
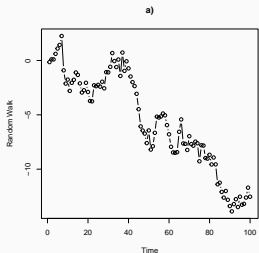
```
> x <- rnorm(1000)
> y <- rnorm(1000)
> plot(x, y)
```



Plotting using base graphics

- The **type** argument changes the type of plotting done
 - "p" draws points
 - "l" draws lines
 - "o" draws lines and points over-plotted
 - "b" draws lines and points
 - "h" draws histogram-like bars
 - "s" draws stepped lines
- **main** control the title of the plot
- **xlab** & **ylab** control axis labels

```
> x <- 1:100
> y <- cumsum(rnorm(100))
> layout(matrix(1:2, ncol = 1))
> plot(x, y, type = "b", main = "a)", xlab = "Time", ylab = "Random Walk")
> plot(x, y, type = "s", main = "b)", xlab = "Time", ylab = "Random Walk")
> layout(1)
```

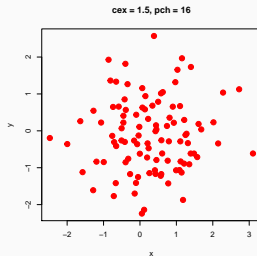
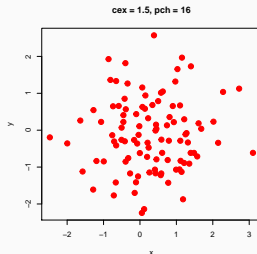


Plotting using base graphics

- `pch` controls the plotting character
- `cex` controls the size of the character
- `col` controls colour
- `axes` logical; should axes be drawn
- `ann` logical; should the plot be annotated
- `axis()`, `title()`, `box()` used to build up plotting
- Allows finer control

```
> x <- rnorm(100)
> y <- rnorm(100)
> layout(matrix(1:2, ncol = 1))
> plot(x, y, main = "cex = 1.5, pch = 16", cex = 1.5, pch = 19, col = "red")
> plot(x, y, cex = 1.5, pch = 19, col = "red", axes = FALSE, ann = FALSE)
> axis(side = 1)
> axis(side = 2)
> title(main = "cex = 1.5, pch = 16", xlab = "x", ylab = "y")
> box()

> layout(1)
```

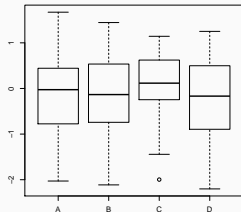
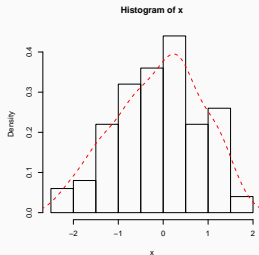


Plotting using base graphics

- `hist()` draws histograms
- `boxplot()` draws boxplots
- `"lwd"` controls the line width
- `"lty"` controls the line type
- `lines()` used to add lines to an existing plot
- Also `points()`

```
> x <- rnorm(100)
> grps <- factor(sample(LETTERS[1:4], 100, replace = TRUE))
> layout(matrix(1:2, ncol = 1))
> dens <- density(x)
> hist(x, freq = FALSE)
> lines(dens, col = "red", lwd = 2, lty = "dashed")
> boxplot(x ~ grps)

> layout(1)
```



Plotting device regions and margins



Plotting device regions and margins

- Control the size of margins using several parameters
 - `mar` — set margins in terms of number of lines of text
 - `mai` — set margins in terms of number of inches
- Specify as a vector of length 4 — `mar = c(5,4,4,2) + 0.1`
- The ordering is Bottom, Left, Top, Right
- The outer margin is controlled via parameter `oma` and `omi`, just like `mar`
- By default, there is no outer margin — `oma = c(0,0,0,0)`

```
> x <- runif(100)
> y <- 4 + (2.1 * x) + rnorm(100, 0, 3)
> op <- par(mar = c(4,4,4,4) + 0.1)
> plot(y ~ x)

> op <- par(op)
>
> x <- runif(100)
> y <- 4 + (2.1 * x) + rnorm(100, 0, 3)
> op <- par(mar = c(4,4,4,4) + 0.1, oma = rep(2,4))
> plot(y ~ x)

> op <- par(op)
```

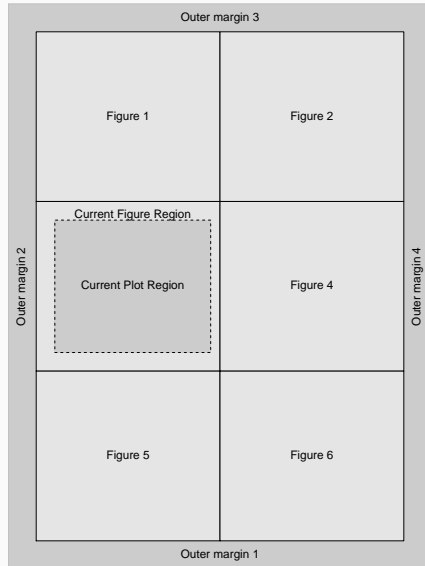
Setting graphical parameters

- Base graphics are controlled by a large number of plotting graphical parameters
- These are detailed in the help page `?par`
- Graphical parameters are changed using the `par()` function and some may be changed within plotting calls
- To avoid getting into a muddle, when changing `par` you should
 - Store the defaults
 - Change your parameters as required
 - When finished the current plot, reset the parameters
- The first two can be done with a single R call

```
> ## Store defaults in 'op' and change current parameters
> op <- par(las = 2, mar = rep(4, 4), oma = c(1,3,4,2), cex.main = 2)
> plot(1:10) # plot something

> par(op)      # reset
```

Plotting on multiple device regions

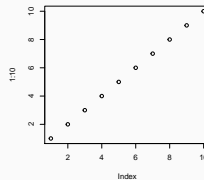
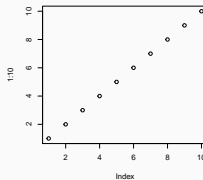
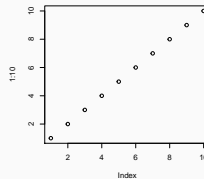
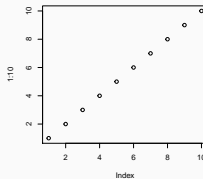


Plotting on multiple device regions

Several ways to split a region into multiple plotting regions

- Graphical parameters `mfrow` & `mfcol`
- The `layout()` function
- The `split.screen()` function

```
> op <- par(mfrow = c(2,2))  
> plot(1:10)  
> plot(1:10)  
> plot(1:10)  
> plot(1:10)  
> par(op)
```



Plotting on multiple device regions

`layout()` allows you to

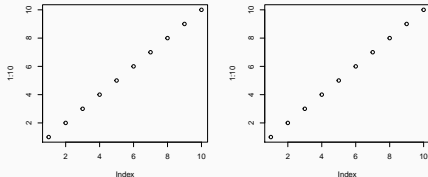
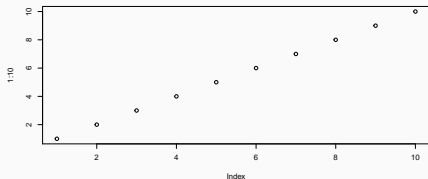
- break the device into multiple regions with individual regions not tied to a single row or column
- specify where each plot should go via an ID

Previous figure could have been produced using

```
layout(matrix(1:4, ncol = 2, byrow = TRUE))
```

Here, we allow the whole first row to be occupied by region 1

```
> layout(matrix(c(1,1,2,3), ncol = 2, byrow = TRUE))  
> plot(1:10)  
> plot(1:10)  
> plot(1:10)  
> layout(1)
```



ggplot

ggplot2 is what all the cool, young kids are using

High-level plotting package like Lattice, but designed for ease of use

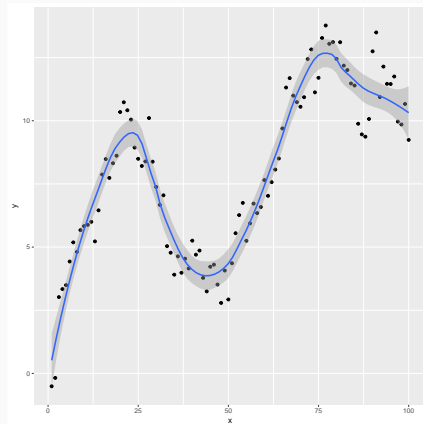
Based on the Leland Wilkinson's *Grammar of Graphics*

Three key principles

- a tidy data structure
- map data to plot features — **aesthetics**
- display data using geometric objects — **geoms**

Plots are built up in layers, composed with +

```
> df <- data.frame(x = 1:100, y = cumsum(rnorm(100)))  
> ggplot(df, aes(x = x, y = y)) +  
+   geom_point() +  
+   geom_smooth(span = 0.3)
```



ggplot — mapping & aesthetics

Variables in the data object are *mapped* to aesthetics that you perceive on the plot

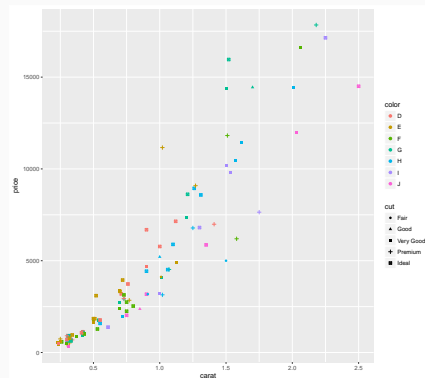
Commonly used aesthetics are

- x and y
- **color** (or **colour**)
- **shape**
- **size**
- **fill** (for anything covering an area, not points)

Specified via the **aes()** function

ggplot handles the creation of legends (keys) for you

```
> data(diamonds)
> set.seed(1410)
> dsmall <- diamonds[sample(nrow(diamonds), 100), ]
> ggplot(dsmall, aes(x = carat, y = price, colour = color, shape = cut))
+   geom_point()
```



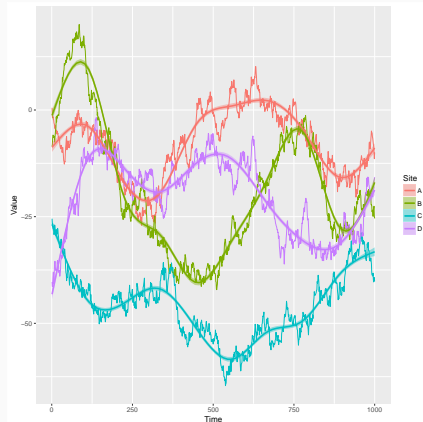
Geometric objects control the way the data are represented on the plot

Common geoms include:

- `geom = "point"` — scatterplot
- `geom = "smooth"` — fits a smooth to the data and draws the smooth and its standard error
- `geom = "boxplot"` — box plots
- `geom = "line"` and `geom = "path"` produce line plots. `"line"` produces lines from left to right, whilst `"path"` draws in data order
- `geom = "histogram"` — histograms
- `geom = "freqpoly"` — frequency polygons
- `geom = "density"` — density plots
- `geom = "bar"` — bar plots

If we have grouping variable(s), we can map them to aesthetics
or we can **facet** the plot to produce multiple panels
Already seen how mapping works
Here we create a random time series for each of four
hypothetical sites

```
> set.seed(789)
> dat2 <- data.frame(x = rep(1:1000, 4),
+                   y = cumsum(rnorm(1000*4)),
+                   Site = factor(rep(LETTERS[1:4], each = 1000)))
> ggplot(dat2, aes(x = x, y = y, colour = Site, fill = Site)) +
+   geom_line() + geom_smooth() +
+   labs(x = "Time", y = "Value")
```



ggplot — facetting

Facetting is what **ggplot** calls the process of breaking up data via one or more grouping variable and displaying a panel for each group that shows the data for that group

Two types of facetting

- `facet_wrap()` — wraps facets into a tabular arrangement
- `facet_grid()` — arranges facets by 1 or 2 categorical variables assigned to the rows and columns of the grid

`facet_wrap()` takes a one-sided formula

`facet_grid()` takes a two-sided formula

```
> ggplot(dat2, aes(x = x, y = y, colour = Site, fill = Site)) +  
+   geom_line() + geom_smooth() +  
+   labs(x = "Time", y = "Value") +  
+   facet_wrap(~ Site, ncol = 2)
```



ggplot — facetting

Facetting is what **ggplot** calls the process of breaking up data via one or more grouping variable and displaying a panel for each group that shows the data for that group

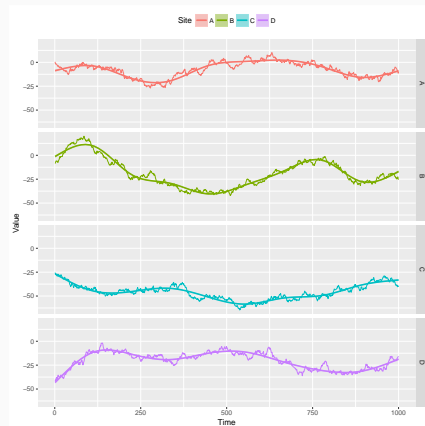
Two types of facetting

- `facet_wrap()` — wraps facets into a tabular arrangement
- `facet_grid()` — arranges facets by 1 or 2 categorical variables assigned to the rows and columns of the grid

`facet_wrap()` takes a one-sided formula

`facet_grid()` takes a two-sided formula

```
> ggplot(dat2, aes(x = x, y = y, colour = Site, fill = Site)) +  
+   geom_line() + geom_smooth() +  
+   labs(x = "Time", y = "Value") +  
+   facet_grid(Site ~ .) +  
+   theme(legend.position = "top")
```



R package management

- CRAN contains hundreds of packages of user-contributed code that you can install from an R session
- Package installation via function `install.packages()`
- Packages can be updated via function `update.packages()`
- When installing or updating for the first time in a session, R will prompt you to choose a mirror to download from
- Once a package is installed you need to load it ready for use
- Load a package from your library using `library()` or `require()`
- Windows and MacOS have menu items to assist with these operations as does RStudio

```
> install.packages("vegan")  
> update.packages()  
> update.packages(ask = FALSE)  
> library("vegan")
```

- It is useful to create your own library for downloaded packages
- This library will not be overwritten when you install a new version of R
- To set a directory you have write permissions on as your user library, create a file named **.Renviron** in your home directory
 - On Windows this is usually `C:\Documents~and~Settings\username\My~Documents`
 - On Linux it is `/home/user/`
- To set your user library to stated directory, add following to your **.Renviron**
 - On Windows if installed R to `C:\R` add: `R_LIBS=C:/R/myRlib`
 - On Linux, create directory `/home/user/R/libs` say and then add: `R_LIBS=/home/user/R/libs`

Copyright © (2015–2017) Gavin L. Simpson Some Rights Reserved

Unless indicated otherwise, this slide deck is licensed under a Creative Commons Attribution 4.0 International License.

