# SPEAKER_NOTES.md - R

Speaker Notes for the 2016-01-12 Software Carpentry R lesson

**TYPE ALL EXAMPLES AS YOU GO. THIS KEEPS THE SPEED SANE, AND ALLOWS YOU TO EXPLAIN EVERY STEP.**

**START SLIDES WITH** `reveal-md slides.md --theme=white`

## `R for reproducible scientific analysis`

**SLIDE** (Learning objectives)

- Welcome
- Teaching

    - Talk around slide

- **Our goal is not just to "do stuff"**

    - do it so that anyone can easily and exactly replicate our workflow and results

# Introduction to R and RStudio

**SLIDE** (Why `R` / `RStudio` ?)

- Talk around slide

**SLIDE** ( `R` / `RStudio` presentation)

- Live presentation section
- Everyone start up `RStudio`

**Summarise windows**

- Four (maybe three) subwindows:
    - Interactive `R` console
    - Editor (may be missing on startup - will appear when files are opened)
    - Environment/History
    - Files/Plots/Packages/Help

## Create a working directory with version control

- **We're following practices of project management**

  - We'll create a project directory, with `Git` version control
  - Helps ensure data integrity
  - Makes sharing code easier (lab-mates, publication)
  - Easier to recover after a Christmas break

- **Create the new directory**

  - `File->New Project`
  - `New Directory`
  - `Empty Project`
  - Enter sensible name, e.g. `swc-workshop`
  - Check box for `Create a git repository`
  - `Create project`

**GREEN/RED STICKY CHECK**

- Describe contents of new folder
  - `.gitignore`
  - `.Rproj`

**SLIDE** (Best practices)

- Talk around slide

# Create directory structure

**SLIDE** (Creating files/directories)

- Live presentation section

- **Create subdirectory for data**

  - In `Files` tab, create `data` subdirectory

- **Create new `R` script**

  - `File -> New File -> R script`
  - save in working directory with sensible name, e.g. `swc-script.R`

**GREEN/RED STICKY CHECK**

# Version control

- **Show Git tab on right**

- **Stage files**

  - Three files shown (including `.gitignore` and the new script file)
  - Yellow status markers mean they're not in the repository
  - Click check-boxes to stage them
  - Note that **we don't version disposable output**

- **Commit files**

  - Click `Commit`
  - Describe new dialogue window
  - Show contents/changes to files
  - Add commit message ("Initialised repository")
  - Commit
  - Show commit summary
  - Exit

**GREEN/RED STICKY CHECK**

**SLIDE** (Challenge 1)

Run through challenge (5min?) - hint about editing `.gitignore`

- Right-click link on presentation and download to `data`
- Create `graphs` subdirectory in `Files` tab
- Edit `.gitignore` to add `graphs/` folder and save
- Stage `.gitignore` in Git tab
- Commit in Git tab, and add appropriate commit message
- Demo History window for Git

**SLIDE** ( `R` as a calculator)

# Interacting with `R`

- **Two ways**

  - Type commands in the console
  - Use the script editor and save the script

- **Console**

  - Output shown here
  - Good for experimentation
  - Commands 'forgotten' when you close a session

- **Script**

- Keeps record of what you did
- Easier to reproduce and share

## Working at the console

- `R` shows a `>` if it is expecting input

```
1  > 1 + 100
2  [1] 101
```

- `R` shows `+` if it's waiting for completion (`Esc` to exit)

```
1  > 1 +
2  +
```

## Working from script file

- Can write same commands in the script file ( `1 + 100` )

  - Use `Run` to execute
  - Use `Ctrl-Enter` to execute
  - Output appears in the console
  - Show `#` comments - good practice to comment
  - More examples (order of precedence):

```
1  > 3 + 5 * 2
2  [1] 13
3  > (3 + 5) * 2
4  [1] 16
```

- Show `Source` operation: runs all script

```
1  > # Using R as a calculator demo
2  > 1 + 100
3  [1] 101
4  > 3 + 5 * 2
5  [1] 13
6  > (3 + 5) * 2
7  [1] 16
```

- More examples
  - scientific notation

```
1  > 1/40
2  [1] 0.025
3  > 2/10000
4  [1] 2e-04
5  > 5e3
6  [1] 5000
```

## Mathematical functions

- General format: `fn(arg)`

    - autocompletion - example: `factorial(6)`

```
1  > sin(1)
2  [1] 0.841471
3  > log(1)
4  [1] 0
5  > log10(10)
6  [1] 1
7  > exp(0.5)
8  [1] 1.648721
```

## Comparisons

- Return `TRUE` / `FALSE` logical values

```
1   > 1 == 1
2   [1] TRUE
3   > 1 == 2
4   [1] FALSE
5   > 1 != 2
6   [1] TRUE
7   > 1 < 2
8   [1] TRUE
9   > 1 > 2
10  [1] FALSE
11  > 1 <= 2
12  [1] TRUE
13  > 1 >= 2
14  [1] FALSE
```

- Computer representation of numbers are approximate: important for comparisons
    - Any physicists/computer scientists in the room?
    - Numbers may not be equal, but be 'the same'
    - Use `all.equal` instead of `==`

```
1  > all.equal(pi-1e-7, pi)
2  [1] "Mean relative difference: 3.183099e-08"
3  > all.equal(pi-1e-8, pi)
4  [1] TRUE
5  > pi-1e-8 == pi
6  [1] FALSE
```

## Variables and assignment

- **Variables hold values**, just like in Python

- Two ways to assign variables

  - The `<-` form is more widely used
  - Consistency more important than choice

```
1  > x <- 1/40
2  > x
3  [1] 0.025
4  > x = 1/40
5  > x
6  [1] 0.025
```

- **Look at the Environment tab** automatic updates

```
1  > x <- 100
```

- Variables can be used as arguments to functions

```
1  > log(x)
2  [1] 4.60517
3  > sqrt(x)
4  [1] 10
```

- Variables can be used to reassign values to themselves

```
1  > x
2  [1] 100
3  > x <- x + 1
4  > x
5  [1] 101
```

**SLIDE** (Good variable names)

- Talk around slide

**SLIDE** (MCQ1)

- Pose question

## Package management

**SLIDE** (Package Management)

- See what packages are installed with `installed.packages()`

    - **demo this one**

- Add a new package using `install.packages("packagename")`

    - **demo this one with** `install.packages("ggplot2")`

- Update packages with `update.packages()`

    - **demo this one**

- You can remove a package with `remove.packages("packagename")`

- To make a package available for use, use `library(packagename)`

    - **demo**

```
1  > ggplot()
2  Error: could not find function "ggplot"
3  > library(ggplot2)
4  Warning message:
5  package 'ggplot2' was built under R version 3.2.3
6  > ggplot()
7  Warning message:
8  In max(vapply(evaled, length, integer(1))) :
9    no non-missing arguments to max; returning -Inf
```

**SLIDE** (Challenge 2)

Solution:

```
1  install.packages("plyr")
2  install.packages("gapminder")
3  install.packages("dplyr")
4  install.packages("tidyr")
```

## Getting help for functions

**SLIDE** (Functions, and getting help)

- Talk around slide

- Demo: `round(3.14159)` :

    - argument: `3.14159`
    - value: `3`

```
1   > round(3.14159)
2   [1] 3
```

**SLIDE** (Getting help for functions)

- **Carrying on with `round()` from last slide**

- What other arguments can `round()` take?

    - Use `args(fname)`

```
1   > args(round)
2   function (x, digits = 0)
3   NULL
```

- Can use the `digits` argument by naming it, or not (but order matters)

```
1   > round(3.14159, digits=2)
2   [1] 3.14
3   > round(3.14159, 2)
4   [1] 3.14
```

- **Best practice**: always use the argument name

    - clearer to others
    - if function changes, order may change
    - difficult to remember the purpose of each argument, if not explicit

- What does a function do?

    - Use `?fname` or `help(fname)` to get the complete help text
    - Demo: `?round` - go through main points

- What package is my function in?

    - (i.e. I can't find it, and don't know what to install)
    - Demo: `??melt` - show that we need `reshape2`

- Is there a function that does X?

    - e.g. you know the name of a test, such as Kolmogorov-Smirnov
    - Demo: `help.search("smirnov")` , `?ks.test`

**SLIDE** (Where can I get more help?)

- Talk around slide

**SLIDE** (Asking the right questions)

- Talk around slide

- For `dput()` example use `dput(head(iris))`

- Demo `sessionInfo()`

# Data Types and Structures in `R`

**SLIDE** (Data Structures in `R` )

- Good place to ask about pace/if a break is needed?

**SLIDE** (Learning Objectives)

- Talk around the slide

- **R is largely used for data analysis**

    - The management and manipulation of data depends on the type of data we have
    - A large amount of day-to-day frustration of learners comes down to problems with data types
    - It's *very important* to understand how `R` sees your data

**SLIDE** (Five "atomic" data types)

- Talk around slide

**SLIDE** (Atomic data types)

- **Create some variables in script**

```
1  # Some variables
2  truth <- TRUE
3  lie <- FALSE
4  i <- 3L
5  d <- 3.0
6  c <- 3 + 0i
7  txt <- "TRUE"
```

- Show equivalence of integer, double and complex

```
1   > typeof(i)
2   [1] "integer"
3   > typeof(d)
4   [1] "double"
5   > i == c
6   [1] TRUE
7   > d == c
8   [1] TRUE
9   > i == d
10  [1] TRUE
11  > is.numeric(i)
12  [1] TRUE
13  > is.numeric(d)
14  [1] TRUE
15  > is.numeric(c)
16  [1] FALSE
```

- Show other types

```
1   > typeof(truth)
2   [1] "logical"
3   > typeof(lie)
4   [1] "logical"
5   > typeof(txt)
6   [1] "character"
```

- `is.X()` tests for a data type

```
1   > is.logical(lie)
2   [1] TRUE
3   > is.logical(txt)
4   [1] FALSE
5   > is.integer(i)
6   [1] TRUE
7   > is.integer(d)
8   [1] FALSE
```

**SLIDE** (Five data structures)

- Talk around slide
  - more on `data.frame` in detail later

# Vectors

**SLIDE** (Vectors)

- Vectors are the most common data structure
- Vectors can contain only one data type

- vectors also known as "atomic vectors"

- **The `c()` function**

  - `c()` is the "concatenate" function

```
1  > x <- c(10, 12, 45, 33)
2  > x
3  [1] 10 12 45 33
```

- **Number sequences**

  - can use `:` or `seq()` functions
  - both functions *return* vectors

```
1  > series <- 1:10
2  > series
3   [1]  1  2  3  4  5  6  7  8  9 10
4  > series <- seq(15)
5  > series
6   [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15
7  > seq(1, 10, by=0.5)
8   [1]  1.0  1.5  2.0  2.5  3.0  3.5  4.0  4.5  5.0  5.5  6.0  6.5  7.0  7.5  8.0  8.5
9  [18]  9.5 10.0
```

- **What type is our vector?**

  - Use the `str()` (structure) function

```
1  > str(x)
2   num [1:4] 10 12 45 33
3  > str(series)
4   int [1:15] 1 2 3 4 5 6 7 8 9 10 ...
5  > is.numeric(x)
6  [1] TRUE
7  > is.numeric(series)
8  [1] TRUE
9  > is.integer(x)
10 [1] FALSE
11 > is.integer(series)
12 [1] TRUE
```

- **Series is `integer` type, but `x` is not**

  - The `c()` function automatically turns integers into 'real'/'double' numbers
  - To specify integers, use `L` :

```
1  > y <- c(10L, 12L, 45L, 33L)
2  > y
3  [1] 10 12 45 33
4  > x
5  [1] 10 12 45 33
6  > is.integer(x)
7  [1] FALSE
8  > is.integer(y)
9  [1] TRUE
```

- **Extending a vector**

  - Append new elements to a vector with `c()`

```
1  > x
2  [1] 10 12 45 33
3  > x <- c(x, 57)
4  > x
5  [1] 10 12 45 33 57
```

- **Character vectors**

  - You can use `c()` to create vectors from any datatype, including characters

```
1  > t <- c('a', 'b', 'c')
2  > t
3  [1] "a" "b" "c"
4  > str(t)
5   chr [1:3] "a" "b" "c"
```

**SLIDE** (Challenge 2)

- Point out that `R` will attempt to "coerce" the datatype to be one that can represent all items in the vector.

Solution:

```
1  > xx <- c(1.7, 'a')
2  > str(xx)
3   chr [1:2] "1.7" "a"
4  > xx <- c(TRUE, 2)
5  > str(xx)
6   num [1:2] 1 2
7  > xx <- c('a', TRUE)
8  > str(xx)
9   chr [1:2] "a" "TRUE"
```

**SLIDE** (Coercion)

- Talk around slide

- **DEMO**

```
1   > x
2   [1] 10 12 45 33 57
3   > str(x)
4    num [1:5] 10 12 45 33 57
5   > as.character(x)
6   [1] "10" "12" "45" "33" "57"
7   > as.complex(x)
8   [1] 10+0i 12+0i 45+0i 33+0i 57+0i
9   > as.logical(x)
10  [1] TRUE TRUE TRUE TRUE TRUE
```

- Sometimes coercion is not possible

```
1   > x <- c('a', 'b', 'c')
2   > str(x)
3    chr [1:3] "a" "b" "c"
4   > as.numeric(x)
5   [1] NA NA NA
6   Warning message:
7   NAs introduced by coercion
```

**SLIDE** (Useful vector functions)

- There are functions that will give information about the vector

```
1   > x <- 0:10
2   > tail(x)
3   [1]  5  6  7  8  9 10
4   > tail(x, n=2)
5   [1]  9 10
6   > head(x)
7   [1] 0 1 2 3 4 5
8   > head(x, n=2)
9   [1] 0 1
10  > length(x)
11  [1] 11
12  > str(x)
13   int [1:11] 0 1 2 3 4 5 6 7 8 9 ...
```

- Vector elements can also be named (this is *similar to*, but not the same as a Python dictionary)

```
1   > x <- 1:4
2   > names(x)
3   NULL
4   > str(x)
5    int [1:4] 1 2 3 4
6   > names(x) <- c('a', 'b', 'c', 'd')
7   > x
8   a b c d
9   1 2 3 4
10  > str(x)
11   Named int [1:4] 1 2 3 4
12   - attr(*, "names")= chr [1:4] "a" "b" "c" "d"
```

# Factors

**SLIDE** (Factors)

- Talk around slide

**SLIDE** (Factors demo)

- **Create factor**

    - Use the `factor()` function with a vector as the argument
    - Predefined values are those present on creation
    - **Typos can give unexpected levels!**

```
1   > x <- factor(c('yes', 'no', 'no', 'yes', 'yes'))
2   > x
3   [1] yes no  no  yes yes
4   Levels: no yes
5   > levels(x)
6   [1] "no"  "yes"
7   > str(x)
8    Factor w/ 2 levels "no","yes": 2 1 1 2 2
```

- **Ordering levels**

    - Level order may be important
    - Models expect the baseline/control to be the first level
    - By default, `factor()` orders factors alphabetically

```
1   > x <- factor(c('case', 'control', 'control', 'case'))
2   > x
3   [1] case    control control case
4   Levels: case control
5   > levels(x)
6   [1] "case"    "control"
```

- Here, `case` will be considered the baseline/control factor.
- This is not what modelling functions expect - results will be difficult to interpret.
    - Use the `levels=` argument to fix

```
1   > x <- factor(c('case', 'control', 'control', 'case'), levels=c('control', 'case'))
2   > x
3   [1] case    control control case
4   Levels: control case
5   > levels(x)
6   [1] "control" "case"
```

- `table()` and `barplot()` functions

    - The `table()` function can be used to tabulate the number of members of each category
    - Introduces the `Plots` tab for output

```
1   > expt <- factor(c('a', 'b', 'a', 'c', 'a', 'control', 'a', 'b', 'c'))
2   > str(expt)
3    Factor w/ 4 levels "a","b","c","control": 1 2 1 3 1 4 1 2 3
4   > table(expt)
5   expt
6         a       b       c control
7         4       2       2       1
8   > barplot(table(expt))
```

# Matrices

**SLIDE** (Matrices)

- **Creating a matrix**

    - Matrices are essentially atomic vectors with extra dimensions
    - `set.seed()` makes our pseudorandom numbers reproducible
    - `rnorm()` selects values from a standard normal distribution
    - Create matrix with `matrix()`, passing a vector and specifying the number of rows and columns

```
1   > set.seed(1)
2   > x <- matrix(rnorm(18), ncol=6, nrow=3)
3   > x
4             [,1]       [,2]      [,3]       [,4]       [,5]        [,6]
5   [1,] -0.6264538  1.5952808 0.4874291 -0.3053884 -0.6212406 -0.04493361
6   [2,]  0.1836433  0.3295078 0.7383247  1.5117812 -2.2146999 -0.01619026
7   [3,] -0.8356286 -0.8204684 0.5757814  0.3898432  1.1249309  0.94383621
8   > str(x)
9    num [1:3, 1:6] -0.626 0.184 -0.836 1.595 0.33 ...
```

- `RStudio` treats vectors as 'Values' and matrices as 'Data', in the environment
- `RStudio` also lets you see the matrix in the editor window (**demo this**)

```
1  > str(x)
2   num [1:3, 1:6] -0.626 0.184 -0.836 1.595 0.33 ...
3  > length(x)
4  [1] 18
5  > nrow(x)
6  [1] 3
7  > ncol(x)
8  [1] 6
```

**SLIDE** (Challenge 3)

Solution:

```
1  > m <- matrix(1:50, ncol=5, nrow=10)
2  > m
3        [,1] [,2] [,3] [,4] [,5]
4   [1,]    1   11   21   31   41
5   [2,]    2   12   22   32   42
6   [3,]    3   13   23   33   43
7   [4,]    4   14   24   34   44
8   [5,]    5   15   25   35   45
9   [6,]    6   16   26   36   46
10  [7,]    7   17   27   37   47
11  [8,]    8   18   28   38   48
12  [9,]    9   19   29   39   49
13  [10,]  10   20   30   40   50
14  > ?matrix
15  > m <- matrix(1:50, ncol=5, nrow=10, byrow=TRUE)
16  > m
17        [,1] [,2] [,3] [,4] [,5]
18  [1,]    1    2    3    4    5
19  [2,]    6    7    8    9   10
20  [3,]   11   12   13   14   15
21  [4,]   16   17   18   19   20
22  [5,]   21   22   23   24   25
23  [6,]   26   27   28   29   30
24  [7,]   31   32   33   34   35
25  [8,]   36   37   38   39   40
26  [9,]   41   42   43   44   45
27  [10,]  46   47   48   49   50
```

# Lists

- **Creating a list**

- Directly with `list()`
- By coercion with `as.list()`
- Elements indicated/recovered by double-brackets: `[[ ]]`
- **Numbering is 1-based - not like Python/other languages**

```
 1  > x <- list(1, 'a', TRUE, 1+4i)
 2  > x
 3  [[1]]
 4  [1] 1
 5  [[2]]
 6  [1] "a"
 7  [[3]]
 8  [1] TRUE
 9  [[4]]
10  [1] 1+4i
11  > x[[3]]
12  [1] TRUE
```

- elements can be named
  - named elements can be recovered with `$` notation

```
 1  > xlist <- list(a="SWC Workshop", b=1:10, data=head(iris))
 2  > xlist
 3  $a
 4  [1] "SWC Workshop"
 5  $b
 6   [1]  1  2  3  4  5  6  7  8  9 10
 7  $data
 8    Sepal.Length Sepal.Width Petal.Length Petal.Width Species
 9  1          5.1         3.5          1.4         0.2  setosa
10  2          4.9         3.0          1.4         0.2  setosa
11  3          4.7         3.2          1.3         0.2  setosa
12  4          4.6         3.1          1.5         0.2  setosa
13  5          5.0         3.6          1.4         0.2  setosa
14  6          5.4         3.9          1.7         0.4  setosa
15  > xlist$a
16  [1] "SWC Workshop"
```

**SLIDE** (Challenge 4)

Solution:

```
1   > my_list <- list(
2   +   data_types=c("logical", "integer", "double", "complex", "character"),
3   +   data_structures=c("vector", "matrix", "factor", "list")
4   + )
5   > str(my_list)
6   List of 2
7    $ data_types     : chr [1:5] "logical" "integer" "double" "complex" ...
8    $ data_structures: chr [1:4] "vector" "matrix" "factor" "list"
```

# Functions

**SLIDE** (Functions)

**SLIDE** (Learning objectives)

- Talk around slide

- **Why functions?**

    - You've already seen the power of functions, for encapsulating complex analyses into simple commands
    - Functions work similarly in `R` to in Python

**SLIDE** (What is a function?)

- Talk around slide

## Defining a function

**SLIDE** (Defining a function)

- Talk around slide

- **Create a new `R` script file to hold functions**

    - `File -> New File -> R Script`
    - `File -> Save -> functions-lesson.R`
    - Check what's happened in Git tab

- **Write new function in script**

    - Describe parts of function:
    - *prototype* with inputs
    - code block/body
    - indentation (readability)
    - addition, and return statements

- function scope, internal variables (readability)
- assignment of function to variable
- comments (readability)

```
1   # Returns sum of two inputs
2   my_sum <- function(a, b) {
3     the_sum <- a + b
4     return(the_sum)
5   }
6   # Converts fahrenheit to Kelvin
7   fahr_to_kelvin <- function(temp) {
8     kelvin <- ((temp - 32) * (5 / 9)) + 273.15
9     return(kelvin)
10  }
```

- **Run the functions**

  - `source` the script
  - tab-completion works!
  - boiling and freezing points

```
1   > fahr_to_kelvin(32)
2   [1] 273.15
3   > fahr_to_kelvin(212)
4   [1] 373.15
```

**SLIDE** (Challenge 1)

Solution:

```
1   kelvin_to_celsius <- function(temp) {
2     celsius <- temp - 273.15
3     return(celsius)
4   }
```

**SLIDE** (Challenge 2)

Solution:

```
1   fahr_to_celsius <- function(temp) {
2     kelvin <- fahr_to_kelvin(temp)
3     celsius <- kelvin_to_celsius(kelvin)
4     return(celsius)
5   }
```

**INSERTED EXAMPLE**

- Just as in Python, we can use `for` loops to apply a function to several values

- Avoids repetition

```
1  for (i in 32:100) {
2    print(fahr_to_celsius(i))
3  }
```

- Can also apply functions to vectors

```
1  fahr_to_celsius(32:100)
```

- Also `if` and `if/else` statements, as in Python:

```
1  if (5 > 1) {
2    print("condition is true")
3  }
```

```
1  if (5 < 1) {
2    print("condition is true")
3  } else {
4    print("condition is false")
5  }
```

- **Commit to local Git repo**

**SLIDE** (Testing functions)

- Talk around slide

- **Known good values**

  - water freezes at 32F/0C, boils at 212F/100C

```
1  > fahr_to_celsius(32)
2  [1] 0
3  > fahr_to_celsius(212)
4  [1] 100
```

- **Known bad values**

  - All values are fair game on Fahrenheit/Celsius, but can't go below 0K

```
1  > kelvin_to_celsius(-10)
2  [1] -283.15
```

- **We'd need to modify this for real use!**

# Data Frames

**SLIDE** (Data Frames)

**SLIDE** (Learning Objectives)

- Talk around slide

**SLIDE** ( `data.frame` s)

- Talk around slide

**SLIDE** (My first `data.frame` )

- **Create a `data.frame`**

  - Name columns explicitly

```
1   > df <- data.frame(id=c('a', 'b', 'c', 'd', 'e', 'f'), x=1:6, y=c(214:219))
2   > df
3     id x   y
4   1  a 1 214
5   2  b 2 215
6   3  c 3 216
7   4  d 4 217
8   5  e 5 218
9   6  f 6 219
10  > length(df)
11  [1] 3
12  > dim(df)
13  [1] 6 3
14  > ncol(df)
15  [1] 3
16  > nrow(df)
17  [1] 6
18  > summary(df)
19   id            x                y
20   a:1   Min.   :1.00   Min.   :214.0
21   b:1   1st Qu.:2.25   1st Qu.:215.2
22   c:1   Median :3.50   Median :216.5
23   d:1   Mean   :3.50   Mean   :216.5
24   e:1   3rd Qu.:4.75   3rd Qu.:217.8
25   f:1   Max.   :6.00   Max.   :219.0
```

- Rows are named automatically, by default.
- The length of a `data.frame` is the number of columns it has
- Use `dim()` , `nrow()` , `ncol()` to get the numbers of rows and columns

- **`data.frame` s coerce strings/characters to become factors!**

- ◦ We don't always want this
- ◦ Can use the `stringsAsFactors` argument to change this behaviour

```
 1   > str(df)
 2   'data.frame':   6 obs. of  3 variables:
 3    $ id: Factor w/ 6 levels "a","b","c","d",..: 1 2 3 4 5 6
 4    $ x : int  1 2 3 4 5 6
 5    $ y : int  214 215 216 217 218 219
 6   > df <- data.frame(id=c('a', 'b', 'c', 'd', 'e', 'f'), x=1:6, y=c(214:219),
 7                      stringsAsFactors=FALSE)
 8   > df
 9     id x   y
10   1  a 1 214
11   2  b 2 215
12   3  c 3 216
13   4  d 4 217
14   5  e 5 218
15   6  f 6 219
16   > str(df)
17   'data.frame':   6 obs. of  3 variables:
18    $ id: chr  "a" "b" "c" "d" ...
19    $ x : int  1 2 3 4 5 6
20    $ y : int  214 215 216 217 218 219
```

- **Show `data.frame` in Editor tab**

**SLIDE** (Challenge 1)

(5min)

- Solution:
  - ◦ missing quotes in `author_last`
  - ◦ missing date in `year`

```
1   author_book <- data.frame(author_first=c("Charles", "Ernst",
2                                             "Theodosius"),
3                       author_last=c("Darwin", "Mayr", "Dobzhansky"),
4                       year=c(1859, 1942, 1970))
```

**SLIDE** (Challenge 2)

(5min)

Solution:

```
1  > str(country_climate)
2  'data.frame':   4 obs. of  5 variables:
3   $ country            : Factor w/ 4 levels "Australia","Canada",..: 2 3 4 1
4   $ climate            : Factor w/ 4 levels "cold","hot","hot/temperate",..: 1 2 4 3
5   $ temperature        : Factor w/ 4 levels "10","15","18",..: 1 4 3 2
6   $ northern_hemisphere: Factor w/ 2 levels "FALSE","TRUE": 2 2 1 1
7   $ has_kangaroo       : num  0 0 0 1
```

## Adding rows and columns

- **Adding a column with `cbind()`**

  - By default the column doesn't get a name
  - to provide a name, use the name as an argument

```
1  > df
2    id x    y
3  1  a 1 214
4  2  b 2 215
5  3  c 3 216
6  4  d 4 217
7  5  e 5 218
8  6  f 6 219
9  > df <- cbind(df, 6:1)
10 > df
11   id x    y 6:1
12 1  a 1 214   6
13 2  b 2 215   5
14 3  c 3 216   4
15 4  d 4 217   3
16 5  e 5 218   2
17 6  f 6 219   1
18 > df <- cbind(df, caps=LETTERS[1:6])
19 > df
20   id x    y 6:1 caps
21 1  a 1 214   6    A
22 2  b 2 215   5    B
23 3  c 3 216   4    C
24 4  d 4 217   3    D
25 5  e 5 218   2    E
26 6  f 6 219   1    F
```

- Note that `caps` is a factor:

```
1  > LETTERS
2   [1] "A" "B" "C" "D" "E" "F" "G" "H" "I" "J" "K" "L" "M" "N" "O" "P" "Q" "R" "S" "T"
3  [23] "W" "X" "Y" "Z"
4  > typeof(LETTERS)
5  [1] "character"
6  > str(df)
7  'data.frame':   6 obs. of  5 variables:
8   $ id  : chr  "a" "b" "c" "d" ...
9   $ x   : int  1 2 3 4 5 6
10  $ y   : int  214 215 216 217 218 219
11  $ 6:1 : int  6 5 4 3 2 1
12  $ caps: Factor w/ 6 levels "A","B","C","D",..: 1 2 3 4 5 6
```

- **Renaming a column**

    o  Use `names()` or `colnames()` to change the name of a column

```
1  > colnames(df)
2  [1] "id"   "x"    "y"    "6:1"  "caps"
3  > colnames(df)[4]
4  [1] "6:1"
5  > colnames(df)[4] <- 'SixToOne'
6  > colnames(df)
7  [1] "id"       "x"        "y"        "SixToOne" "caps"
```

- **Adding a row with** `rbind`

    o  Add a list (multiple types across columns)
    o  Need to take care that datatypes match the columns of the `data.frame`
    o  *Particularly a problem with characters and factors!*

```
1  > df <- rbind(df, list('g', 11, 42, 0, 'G'))
2  Warning message:
3  In `[<-.factor`(`*tmp*`, ri, value = "G") :
4    invalid factor level, NA generated
5  > df
6    id  x   y SixToOne caps
7  1  a  1 214        6    A
8  2  b  2 215        5    B
9  3  c  3 216        4    C
10 4  d  4 217        3    D
11 5  e  5 218        2    E
12 6  f  6 219        1    F
13 7  g 11  42        0 <NA>
```

- `R` *tried to be helpful, and put a* `NA` *special value to indicate missing data*
- Two options to add the data:

- Coerce the column to be a `character` type
- Add a level to the factor in that column (mostly what we want to do)

```
1  > str(df$caps)
2   Factor w/ 6 levels "A","B","C","D",..: 1 2 3 4 5 6 NA
3  > levels(df$caps)
4  [1] "A" "B" "C" "D" "E" "F"
5  > c(levels(df$caps), 'G')
6  [1] "A" "B" "C" "D" "E" "F" "G"
7  > levels(df$caps) <- c(levels(df$caps), 'G')
8  > str(df$caps)
9   Factor w/ 7 levels "A","B","C","D",..: 1 2 3 4 5 6 NA
```

- Now we can add the row

```
1  > df <- rbind(df, list('g', 11, 42, 0, 'G'))
2  > df
3    id  x   y SixToOne caps
4  1  a  1 214        6    A
5  2  b  2 215        5    B
6  3  c  3 216        4    C
7  4  d  4 217        3    D
8  5  e  5 218        2    E
9  6  f  6 219        1    F
10 7  g 11  42        0 <NA>
11 8  g 11  42        0    G
```

- But we have a problem:
  - There's an `<NA>` in the data that we don't want
  - This can happen in many different ways for real data
  - We'll deal with this in the next section

## Reading in data

**SLIDE** (Reading in data)

- **Most of the time you work with pre-prepared data**

  - We don't often have to build `data.frame`s by hand
  - Most data likely to come from software in a standard form
  - Sometimes it's not in the best condition, though…

- **Inspecting data in file**

  - `Files` tab: navigate to data file
  - click on file

- **Discuss data**

  - Point out comma-separations (not always best choice - Euro data)
  - Point out header line
  - Inspecting the structure of the data means we can specify proper arguments in `read.table`

- **Read data**

  - Using `read.table`
  - Put in script

```
# Load gapminder data
gapminder <- read.table(
  file="data/gapminder-FiveYearData.csv",
  header=TRUE, sep=","
)
head(gapminder)
      country year      pop continent lifeExp gdpPercap
1 Afghanistan 1952  8425333      Asia  28.801  779.4453
2 Afghanistan 1957  9240934      Asia  30.332  820.8530
3 Afghanistan 1962 10267083      Asia  31.997  853.1007
4 Afghanistan 1967 11537966      Asia  34.020  836.1971
5 Afghanistan 1972 13079460      Asia  36.088  739.9811
6 Afghanistan 1977 14880372      Asia  38.438  786.1134
> str(gapminder)
'data.frame':   1704 obs. of  6 variables:
 $ country  : Factor w/ 142 levels "Afghanistan",..: 1 1 1 1 1 1 1 1 1 1 ...
 $ year     : int  1952 1957 1962 1967 1972 1977 1982 1987 1992 1997 ...
 $ pop      : num  8425333 9240934 10267083 11537966 13079460 ...
 $ continent: Factor w/ 5 levels "Africa","Americas",..: 3 3 3 3 3 3 3 3 3 3 ...
 $ lifeExp  : num  28.8 30.3 32 34 36.1 ...
 $ gdpPercap: num  779 821 853 836 740 ...
```

- **Load a dataset from the internet**

  - Using `read.csv` - a special case of `read.table`
  - Automatically uses `sep=","` and `header=TRUE` (not all data well-behaved!)
  - Files need not be local - can use a URL for online data
  - Put in script

```
 1   > # Load survey data
 2   > surveys <- read.csv('http://files.figshare.com/2236372/combined.csv')
 3   > head(surveys)
 4     record_id month day year plot_id species_id sex hindfoot_length weight    genus   spe
 5   1         1     7  16 1977       2         NL   M              32     NA Neotoma albi
 6   2        72     8  19 1977       2         NL   M              31     NA Neotoma albi
 7   3       224     9  13 1977       2         NL                  NA     NA Neotoma albi
 8   4       266    10  16 1977       2         NL                  NA     NA Neotoma albi
 9   5       349    11  12 1977       2         NL                  NA     NA Neotoma albi
10   6       363    11  12 1977       2         NL                  NA     NA Neotoma albi
11       taxa plot_type
12   1 Rodent   Control
13   2 Rodent   Control
14   3 Rodent   Control
15   4 Rodent   Control
16   5 Rodent   Control
17   6 Rodent   Control
18   > str(surveys)
19   'data.frame':    34786 obs. of  13 variables:
20    $ record_id      : int  1 72 224 266 349 363 435 506 588 661 ...
21    $ month          : int  7 8 9 10 11 11 12 1 2 3 ...
22    $ day            : int  16 19 13 16 12 12 10 8 18 11 ...
23    $ year           : int  1977 1977 1977 1977 1977 1977 1977 1978 1978 1978 ...
24    $ plot_id        : int  2 2 2 2 2 2 2 2 2 2 ...
25    $ species_id     : Factor w/ 48 levels "AB","AH","AS",..: 16 16 16 16 16 16 16 16 16
26    $ sex            : Factor w/ 6 levels "","F","M","P",..: 3 3 1 1 1 1 1 1 3 1 ...
27    $ hindfoot_length: int  32 31 NA NA NA NA NA NA NA NA ...
28    $ weight         : int  NA NA NA NA NA NA NA NA 218 NA ...
29    $ genus          : Factor w/ 26 levels "Ammodramus","Ammospermophilus",..: 13 13 13
30    $ species        : Factor w/ 40 levels "albigula","audubonii",..: 1 1 1 1 1 1 1 1 1
31    $ taxa           : Factor w/ 4 levels "Bird","Rabbit",..: 4 4 4 4 4 4 4 4 4 4 ...
32    $ plot_type      : Factor w/ 5 levels "Control","Long-term Krat Exclosure",..: 1 1 1
```

- **Good point to revisit staging/committing to local repo**

    - Go to Git tab
    - Stage current script
    - Inspect with Diff - see what's changed
    - Add commit message
    - Commit

# Indexing and Subsetting data

**SLIDE** (Indexing and Subsetting data)

**SLIDE** (Learning outcomes)

- **We don't always need to use all of the data**

  - There might be incomplete or inappropriate data we need to skip
  - We may only care about a subset of samples/observations
  - We typically want to run cross-validation of statistical models

- Talk around slide

# Subset by index

**SLIDE** (Subset by index)

- **Every element in a collection is indexed**

  - Each item in a collection can be referred to by the index
  - Demonstrate with a vector:

```
1  > x <- c(5.4, 6.2, 7.1, 4.8, 7.5)
2  > names(x) <- letters[1:5]
3  > x
4    a   b   c   d   e
5  5.4 6.2 7.1 4.8 7.5
```

- We can extract elements
  - individually
  - in groups
  - as a 'slice'

- **NOTE: Elements are numbered from 1, not 0 (unlike Python)**

```
1   > x
2     a   b   c   d   e
3   5.4 6.2 7.1 4.8 7.5
4   > x[1]
5     a
6   5.4
7   > x[4]
8     d
9   4.8
10  > x[c(2,4)]
11    b   d
12  6.2 4.8
13  > x[1:3]
14    a   b   c
15  5.4 6.2 7.1
16  > x[c(1,1,3)]
17    a   a   c
18  5.4 5.4 7.1
```

- **Asking for an element that isn't there**

  - `x[0]` gives an empty vector
  - `x[6]` gives a missing value `NA`

```
1   > x[0]
2   named numeric(0)
3   > x[6]
4   <NA>
5     NA
```

## Skip/remove by index

- **Use a negative number to return all *other* elements**

```
1   > x
2     a   b   c   d   e
3   5.4 6.2 7.1 4.8 7.5
4   > x[-2]
5     a   c   d   e
6   5.4 7.1 4.8 7.5
7   > x[c(-1,-5)]
8     b   c   d
9   6.2 7.1 4.8
```

- **Assign the result back to the original collection to *remove* elements**

```
1  > x
2    a   b   c   d   e
3  5.4 6.2 7.1 4.8 7.5
4  > x <- x[-4]
5  > x
6    a   b   c   e
7  5.4 6.2 7.1 7.5
```

**SLIDE** (Challenge 1)

Solution:

```
1   > x[-1:3]
2   Error in x[-1:3] : only 0's may be mixed with negative subscripts
3   > -1:3
4   [1] -1  0  1  2  3
5   > 1:3
6   [1] 1 2 3
7   > -(1:3)
8   [1] -1 -2 -3
9   > x[-(1:3)]
10    d   e
11  4.8 7.5
```

## Logical masks

**SLIDE** (Logical masks)

- Talk around slide

- **Logical mask vectors**

    - Any vector of `TRUE` / `FALSE` values the same size as the vector we subset works
    - Shorter vectors cycle round

```
1   > x <- c(5.4, 6.2, 7.1, 4.8, 7.5)
2   > names(x) <- letters[1:5]
3   > x
4     a   b   c   d   e
5   5.4 6.2 7.1 4.8 7.5
6   > mask <- c(TRUE, FALSE, TRUE, FALSE, TRUE)
7   > mask
8   [1]  TRUE FALSE  TRUE FALSE  TRUE
9   > x[mask]
10    a   c   e
11  5.4 7.1 7.5
12  > mask_short = c(FALSE, TRUE)
13  > x[mask_short]
14    b   d
15  6.2 4.8
```

- Any function that generates a logical output can produce a *mask*

  - Can combine comparators with `&` , `|` , `!`

```
1   > x > 7
2       a     b     c     d     e
3   FALSE FALSE  TRUE FALSE  TRUE
4   > x[x > 7]
5     c   e
6   7.1 7.5
7   > (x > 5) & (x < 7)
8       a     b     c     d     e
9    TRUE  TRUE FALSE  TRUE FALSE
10  > x[(x > 5) & (x < 7)]
11    a   b
12  5.4 6.2
```

**SLIDE** (Challenge 2)

Solution:

```
1   (x < 5) | (x > 7)
2       a     b     c     d     e
3   FALSE FALSE  TRUE  TRUE  TRUE
4   > x[(x < 5) | (x > 7)]
5     c   d   e
6   7.1 4.8 7.5
```

# Subset by name

**SLIDE** (Subset by name)

- **Extracting subsets from vectors by name**

    - Can use names directly
    - Can use vectors of names
    - Can't easily skip/remove, this way

```
1  > x['a']
2    a
3  5.4
4  > x[c('b', 'e')]
5    b   e
6  6.2 7.5
7  > x[-c('b', 'e')]
8  Error in -c("b", "e") : invalid argument to unary operator
```

- **Can use logical comparisons**

    - `names() ==` gives a logical vector
    - `names() %in%` for multiple selections

```
1   > names(x)
2   [1] "a" "b" "c" "d" "e"
3   > names(x) == 'c'
4   [1] FALSE FALSE  TRUE FALSE FALSE
5   > x[names(x) == 'c']
6     c
7   7.1
8   > x[names(x) == c('a', 'e')]
9     a
10  5.4
11  Warning message:
12  In names(x) == c("a", "e") :
13    longer object length is not a multiple of shorter object length
14  > names(x) %in% c('a', 'e')
15  [1]  TRUE FALSE FALSE FALSE  TRUE
16  > x[names(x) %in% c('a', 'e')]
17    a   e
18  5.4 7.5
19  > x[!(names(x) %in% c('a', 'c'))]
20    b   d   e
21  6.2 4.8 7.5
```

- **Can use indexing**

    - `which(names())` returns a vector of indexes
    - `==` and `%in%` as before

```
1  > names(x)
2  [1] "a" "b" "c" "d" "e"
3  > names(x) == 'c'
4  [1] FALSE FALSE  TRUE FALSE FALSE
5  > which(names(x) == 'c')
6  [1] 3
7  > x[which(names(x) == 'c')]
8    c
9  7.1
10 > x[which(names(x) %in% c('a', 'c'))]
11   a   c
12 5.4 7.1
13 > x[-which(names(x) %in% c('a', 'c'))]
14   b   d   e
15 6.2 4.8 7.5
```

**SLIDE** (Challenge 3)

(5min)

- Can't use `x['a']` as it only returns a single value

Solution:

```
1  x[names(x) == 'a']
```

# Subsets of matrices

**SLIDE** (Subsets of matrices)

- Talk around slide

- **Create matrix**

```
1  > set.seed(1)
2  > m <- matrix(rnorm(6*4), ncol=4, nrow=6)
3  > m
4            [,1]       [,2]        [,3]        [,4]
5  [1,] -0.6264538  0.4874291 -0.62124058  0.82122120
6  [2,]  0.1836433  0.7383247 -2.21469989  0.59390132
7  [3,] -0.8356286  0.5757814  1.12493092  0.91897737
8  [4,]  1.5952808 -0.3053884 -0.04493361  0.78213630
9  [5,]  0.3295078  1.5117812 -0.01619026  0.07456498
10 [6,] -0.8204684  0.3898432  0.94383621 -1.98935170
```

- **Specify row and column to extract submatrices**

    - can use ranges or subset data

- Does not return data with same indexes!
- Leave a row or column argument blank to retrieve all rows or columns
- Extracting a single row or column returns a vector
- `R` throws an error if indexes are out of bounds

```
> m[3:4, c(3,1)]
           [,1]       [,2]
[1,]  1.12493092 -0.8356286
[2,] -0.04493361  1.5952808
> m[, c(3,1)]
           [,1]       [,2]
[1,] -0.62124058 -0.6264538
[2,] -2.21469989  0.1836433
[3,]  1.12493092 -0.8356286
[4,] -0.04493361  1.5952808
[5,] -0.01619026  0.3295078
[6,]  0.94383621 -0.8204684
> m[3:4,]
           [,1]       [,2]       [,3]      [,4]
[1,] -0.8356286  0.5757814  1.12493092 0.9189774
[2,]  1.5952808 -0.3053884 -0.04493361 0.7821363
> m[,]
           [,1]       [,2]       [,3]       [,4]
[1,] -0.6264538  0.4874291 -0.62124058  0.82122120
[2,]  0.1836433  0.7383247 -2.21469989  0.59390132
[3,] -0.8356286  0.5757814  1.12493092  0.91897737
[4,]  1.5952808 -0.3053884 -0.04493361  0.78213630
[5,]  0.3295078  1.5117812 -0.01619026  0.07456498
[6,] -0.8204684  0.3898432  0.94383621 -1.98935170
> str(m[3:4,])
 num [1:2, 1:4] -0.836 1.595 0.576 -0.305 1.125 ...
> str(m[3,])
 num [1:4] -0.836 0.576 1.125 0.919
 > m[, c(3,6)]
Error in m[, c(3, 6)] : subscript out of bounds
```

# Subsets of lists

**Slide** (Subsets of lists)

- Talk around slide

- **Create list**

  - Inspect content

```
1   > xlist <- list(a="SWC", b=1:10, data=head(iris))
2   > str(xlist)
3   List of 3
4    $ a   : chr "SWC"
5    $ b   : int [1:10] 1 2 3 4 5 6 7 8 9 10
6    $ data:'data.frame':   6 obs. of  5 variables:
7     ..$ Sepal.Length: num [1:6] 5.1 4.9 4.7 4.6 5 5.4
8     ..$ Sepal.Width : num [1:6] 3.5 3 3.2 3.1 3.6 3.9
9     ..$ Petal.Length: num [1:6] 1.4 1.4 1.3 1.5 1.4 1.7
10    ..$ Petal.Width : num [1:6] 0.2 0.2 0.2 0.2 0.2 0.4
11    ..$ Species     : Factor w/ 3 levels "setosa","versicolor",..: 1 1 1 1 1 1
```

- **Extract list**

  - Uses `[` operator
  - essentially slicing
  - returns a list

```
1   > xlist[1]
2   $a
3   [1] "SWC"
4   > xlist[1:2]
5   $a
6   [1] "SWC"
7   $b
8    [1]  1  2  3  4  5  6  7  8  9 10
```

- **Extract element**

  - Uses `[[` operator
  - returns the atomic data type
  - *you can only extract one element at a time*
  - can use the element name

```
1   > xlist[[1]]
2   [1] "SWC"
3   > xlist[[2]]
4    [1]  1  2  3  4  5  6  7  8  9 10
5   > xlist[[1:2]]
6   Error in xlist[[1:2]] : subscript out of bounds
7   > xlist[['data']]
8     Sepal.Length Sepal.Width Petal.Length Petal.Width Species
9   1          5.1         3.5          1.4         0.2  setosa
10  2          4.9         3.0          1.4         0.2  setosa
11  3          4.7         3.2          1.3         0.2  setosa
12  4          4.6         3.1          1.5         0.2  setosa
13  5          5.0         3.6          1.4         0.2  setosa
14  6          5.4         3.9          1.7         0.4  setosa
```

- **Extract by name**

  - Uses the `$` operator (or `[[]]` as above)

```
1  > xlist$a
2  [1] "SWC"
```

- **Extract element contents**

  - Can subset each of the elements in the list, in the same command

```
1  > xlist$data[4,]
2    Sepal.Length Sepal.Width Petal.Length Petal.Width Species
3  4          4.6         3.1          1.5         0.2  setosa
```

# Subsets of `data.frame` s

**SLIDE** (Subsets of `data.frame` s)

- Talk around slide

- **Extract column as dataframe**

  - Use the `[]` operator - returns a dataframe

```
1   > str(gapminder)
2   'data.frame':    1704 obs. of  6 variables:
3    $ country  : Factor w/ 142 levels "Afghanistan",..: 1 1 1 1 1 1 1 1 1 1 ...
4    $ year     : int  1952 1957 1962 1967 1972 1977 1982 1987 1992 1997 ...
5    $ pop      : num  8425333 9240934 10267083 11537966 13079460 ...
6    $ continent: Factor w/ 5 levels "Africa","Americas",..: 3 3 3 3 3 3 3 3 3 3 ...
7    $ lifeExp  : num  28.8 30.3 32 34 36.1 ...
8    $ gdpPercap: num  779 821 853 836 740 ...
9   > head(gapminder[3])
10        pop
11  1  8425333
12  2  9240934
13  3 10267083
14  4 11537966
15  5 13079460
16  6 14880372
17  > head(gapminder["pop"])
18        pop
19  1  8425333
20  2  9240934
21  3 10267083
22  4 11537966
23  5 13079460
24  6 14880372
```

- **Extract column as atomic vector**

  - Use the `[[]]` or `$` operators

```
1   > head(gapminder[[5]])
2   [1] 28.801 30.332 31.997 34.020 36.088 38.438
3   > head(gapminder[["lifeExp"]])
4   [1] 28.801 30.332 31.997 34.020 36.088 38.438
5   > head(gapminder$lifeExp)
6   [1] 28.801 30.332 31.997 34.020 36.088 38.438
```

- **Extract row/column as dataframe**

  - Use two arguments, as for matrices
  - Returns a dataframe if elements are mixed types
  - To get a column dataframe, use `drop=False` argument

```
 1  > gapminder[1:3,]
 2        country year      pop continent lifeExp gdpPercap
 3  1 Afghanistan 1952  8425333     Asia  28.801   779.4453
 4  2 Afghanistan 1957  9240934     Asia  30.332   820.8530
 5  3 Afghanistan 1962 10267083     Asia  31.997   853.1007
 6  > gapminder[3,]
 7        country year      pop continent lifeExp gdpPercap
 8  3 Afghanistan 1962 10267083     Asia  31.997   853.1007
 9  > head(gapminder[, 3, drop=FALSE])
10         pop
11  1  8425333
12  2  9240934
13  3 10267083
14  4 11537966
15  5 13079460
16  6 14880372
```

**SLIDE** (Challenge 4)

(10min)

Solution:

```
 1  > head(gapminder[gapminder$year == 1957,])
 2        country year      pop continent lifeExp gdpPercap
 3  2  Afghanistan 1957  9240934     Asia  30.332    820.853
 4  14      Albania 1957  1476505   Europe  59.280   1942.284
 5  26      Algeria 1957 10270856   Africa  45.685   3013.976
 6  38       Angola 1957  4561361   Africa  31.999   3827.940
 7  50    Argentina 1957 19610538  Americas  64.399   6856.856
 8  62    Australia 1957  9712569   Oceania  70.330  10949.650
 9  > head(gapminder[, -c(1:4)])
10    lifeExp gdpPercap
11  1  28.801   779.4453
12  2  30.332   820.8530
13  3  31.997   853.1007
14  4  34.020   836.1971
15  5  36.088   739.9811
16  6  38.438   786.1134
17  > head(gapminder[gapminder$year %in% c(2002, 2007),])
18        country year      pop continent lifeExp gdpPercap
19  11 Afghanistan 2002 25268405     Asia  42.129   726.7341
20  12 Afghanistan 2007 31889923     Asia  43.828   974.5803
21  23      Albania 2002  3508512   Europe  75.651 4604.2117
22  24      Albania 2007  3600523   Europe  76.423 5937.0295
23  35      Algeria 2002 31287142   Africa  70.994 5288.0404
24  36      Algeria 2007 33333216   Africa  72.301 6223.3675
```

# `data.frame` manipulation with `dplyr`

**SLIDE** ( `data.frame` manipulation with `dplyr` )

**SLIDE** (Learning objectives)

- Talk around slide

**SLIDE** (What and why is `dplyr` ?)

- Talk around slide

**SLIDE** (Split-Apply-Combine)

- Talk around slide

- **A general technique for reducing the amount of repetition in code**

  - good when datasets can be grouped

**SLIDE** (What and why is `dplyr` ?)

- Talk around slide

- **Load** `dplyr`

```
 1   > library(dplyr)
 2
 3   Attaching package: 'dplyr'
 4
 5   The following objects are masked from 'package:stats':
 6
 7       filter, lag
 8
 9   The following objects are masked from 'package:base':
10
11       intersect, setdiff, setequal, union
```

## `select()` and `filter()`

**SLIDE** ( `select()` )

- Talk around figure

**SLIDE** ( `select()` and `filter()` )

- `select()` **keeps only the selected variables/columns**

- Note that we don't use strings for the column names

```
1   > head(gapminder)
2        country year       pop continent lifeExp gdpPercap
3   1 Afghanistan 1952  8425333      Asia  28.801   779.4453
4   2 Afghanistan 1957  9240934      Asia  30.332   820.8530
5   3 Afghanistan 1962 10267083      Asia  31.997   853.1007
6   4 Afghanistan 1967 11537966      Asia  34.020   836.1971
7   5 Afghanistan 1972 13079460      Asia  36.088   739.9811
8   6 Afghanistan 1977 14880372      Asia  38.438   786.1134
9   > head(select(gapminder, year, country, gdpPercap))
10    year     country gdpPercap
11  1 1952 Afghanistan  779.4453
12  2 1957 Afghanistan  820.8530
13  3 1962 Afghanistan  853.1007
14  4 1967 Afghanistan  836.1971
15  5 1972 Afghanistan  739.9811
16  6 1977 Afghanistan  786.1134
```

- **Using the %>% pipe**

  - Analogous to the `|` pipe in the shell
  - Can perform selections without specifying the `data.frame` in the function itself
  - (this *is* useful…)
  - **NOTE: Pipes let us split commands over several lines**

```
1   > year_country_gdp <- gapminder %>% select(year, country, gdpPercap)
2   > head(year_country_gdp)
3    year     country gdpPercap
4   1 1952 Afghanistan  779.4453
5   2 1957 Afghanistan  820.8530
6   3 1962 Afghanistan  853.1007
7   4 1967 Afghanistan  836.1971
8   5 1972 Afghanistan  739.9811
9   6 1977 Afghanistan  786.1134
```

- **Using `filter()` to keep only some data values**

  - Filter lets us restrict rows on the basis of data content

```
1   > head(filter(gapminder, continent=="Europe"))
2     country year     pop continent lifeExp gdpPercap
3   1 Albania 1952 1282697    Europe   55.23  1601.056
4   2 Albania 1957 1476505    Europe   59.28  1942.284
5   3 Albania 1962 1728137    Europe   64.82  2312.889
6   4 Albania 1967 1984060    Europe   66.22  2760.197
7   5 Albania 1972 2263554    Europe   67.69  3313.422
8   6 Albania 1977 2509048    Europe   68.93  3533.004
```

- **Combining** `filter()` **and** `select()` **with pipes**

  - `dplyr` makes combining selection/filtering easy, using pipes
  - Note: we don't need to define an intermediate `data.frame`
  - Note: we don't need to use clunky indexing/names

```
1   > year_country_gdp_euro <- gapminder %>% filter(continent=="Europe")
2   >        %>% select(year, country, gdpPercap)
3   > head(year_country_gdp_euro)
4     year country gdpPercap
5   1 1952 Albania  1601.056
6   2 1957 Albania  1942.284
7   3 1962 Albania  2312.889
8   4 1967 Albania  2760.197
9   5 1972 Albania  3313.422
10  6 1977 Albania  3533.004
```

**SLIDE** (Challenge 1)

Solution:

```
1   > head(gapminder %>% filter(continent=="Africa") %>% select(lifeExp, country, year))
2     lifeExp country year
3   1  43.077 Algeria 1952
4   2  45.685 Algeria 1957
5   3  48.303 Algeria 1962
6   4  51.407 Algeria 1967
7   5  54.518 Algeria 1972
8   6  58.014 Algeria 1977
```

## `group_by()` **and** `summarize`

**SLIDE** (Reducing repetition)

- Talk around slide

**SLIDE** ( `group_by()` )

- Talk round figure
  - separates out `data.frame` on the basis of values in `a`

**SLIDE** ( `summarize()` )

- Talk round figure
  - Creates new variables that repeat over a series of `data.frame` s

**SLIDE** ( `group_by()` and `summarize()` )

- Talk around slide

- `group_by()` produces a "grouped `data.frame`"

  - Not the same as a `data.frame`!
  - Like a `list` where each item is a `data.frame` whose rows correspond only to a particular value of `continent`
  - `tally()` counts up the rows in each group

```
 1  > gapminder %>% group_by(continent)
 2  Source: local data frame [1,704 x 6]
 3  Groups: continent [5]
 4          country  year       pop continent lifeExp gdpPercap
 5           (fctr) (int)     (dbl)    (fctr)   (dbl)     (dbl)
 6  1  Afghanistan  1952   8425333      Asia  28.801  779.4453
 7  2  Afghanistan  1957   9240934      Asia  30.332  820.8530
 8  3  Afghanistan  1962  10267083      Asia  31.997  853.1007
 9  4  Afghanistan  1967  11537966      Asia  34.020  836.1971
10  5  Afghanistan  1972  13079460      Asia  36.088  739.9811
11  6  Afghanistan  1977  14880372      Asia  38.438  786.1134
12  7  Afghanistan  1982  12881816      Asia  39.854  978.0114
13  8  Afghanistan  1987  13867957      Asia  40.822  852.3959
14  9  Afghanistan  1992  16317921      Asia  41.674  649.3414
15  10 Afghanistan  1997  22227415      Asia  41.763  635.3414
16  ..          ...   ...       ...       ...     ...       ...
17  > str(gapminder)
18  'data.frame':    1704 obs. of  6 variables:
19   $ country  : Factor w/ 142 levels "Afghanistan",..: 1 1 1 1 1 1 1 1 1 1 ...
20   $ year     : int  1952 1957 1962 1967 1972 1977 1982 1987 1992 1997 ...
21   $ pop      : num  8425333 9240934 10267083 11537966 13079460 ...
22   $ continent: Factor w/ 5 levels "Africa","Americas",..: 3 3 3 3 3 3 3 3 3 3 ...
23   $ lifeExp  : num  28.8 30.3 32 34 36.1 ...
24   $ gdpPercap: num  779 821 853 836 740 ...
25  > str(gapminder %>% group_by(continent))
26  Classes 'grouped_df', 'tbl_df', 'tbl' and 'data.frame': 1704 obs. of  6 variables:
27   $ country  : Factor w/ 142 levels "Afghanistan",..: 1 1 1 1 1 1 1 1 1 1 ...
28   $ year     : int  1952 1957 1962 1967 1972 1977 1982 1987 1992 1997 ...
29   $ pop      : num  8425333 9240934 10267083 11537966 13079460 ...
30   $ continent: Factor w/ 5 levels "Africa","Americas",..: 3 3 3 3 3 3 3 3 3 3 ...
31   $ lifeExp  : num  28.8 30.3 32 34 36.1 ...
32   $ gdpPercap: num  779 821 853 836 740 ...
33   - attr(*, "vars")=List of 1
34    ..$ : symbol continent
35   - attr(*, "drop")= logi TRUE
36   - attr(*, "indices")=List of 5
37    ..$ : int  24 25 26 27 28 29 30 31 32 33 ...
38    ..$ : int  48 49 50 51 52 53 54 55 56 57 ...
39    ..$ : int  0 1 2 3 4 5 6 7 8 9 ...
```

```
40   ..$ : int   12 13 14 15 16 17 18 19 20 21 ...
41   ..$ : int   60 61 62 63 64 65 66 67 68 69 ...
42  - attr(*, "group_sizes")= int   624 300 396 360 24
43  - attr(*, "biggest_group_size")= int 624
44  - attr(*, "labels")='data.frame':   5 obs. of  1 variable:
45   ..$ continent: Factor w/ 5 levels "Africa","Americas",..: 1 2 3 4 5
46   ..- attr(*, "vars")=List of 1
47   .. ..$ : symbol continent
48   ..- attr(*, "drop")= logi TRUE
49  > gapminder %>%
50  + group_by(continent) %>%
51  + tally()
52  Source: local data frame [5 x 2]
53    continent       n
54       (fctr) (int)
55  1    Africa    624
56  2  Americas    300
57  3      Asia    396
58  4    Europe    360
59  5   Oceania     24
```

- `summarize()` creates summary information for each group

  - We need to tell `summarize()` a function to apply to each of our grouped `data.frame` s
  - We also tell it a variable name to place that calculated value into
  - `summarize` returns a `data.frame`

```
 1   > gapminder %>% group_by(continent)
 2                %>% summarize(meangdpPercap=mean(gdpPercap))
 3   Source: local data frame [5 x 2]
 4
 5     continent meangdpPercap
 6        (fctr)         (dbl)
 7   1    Africa      2193.755
 8   2  Americas      7136.110
 9   3      Asia      7902.150
10   4    Europe     14469.476
11   5   Oceania     18621.609
12   > gapminder %>% group_by(continent)
13                %>% summarize(sdgdpPercap=sd(gdpPercap))
14   Source: local data frame [5 x 2]
15
16     continent sdgdpPercap
17        (fctr)        (dbl)
18   1    Africa     2827.930
19   2  Americas     6396.764
20   3      Asia    14045.373
21   4    Europe     9355.213
22   5   Oceania     6358.983
23   > str(gapminder %>% group_by(continent) %>% summarize(sdgdpPercap=sd(gdpPercap)))
24   Classes 'tbl_df', 'tbl' and 'data.frame':   5 obs. of  2 variables:
25    $ continent  : Factor w/ 5 levels "Africa","Americas",..: 1 2 3 4 5
26    $ sdgdpPercap: num  2828 6397 14045 9355 6359
```

**SLIDE** (Challenge 2)

- Use `filter()` to get the rows you need

Solution:

```
 1   > lifeExp_bycountry <- gapminder %>% group_by(country)
 2                       %>% summarize(meanlifeExp=mean(lifeExp))
 3   > head(lifeExp_bycountry)
 4   Source: local data frame [6 x 2]
 5          country meanlifeExp
 6           (fctr)        (dbl)
 7   1 Afghanistan     37.47883
 8   2      Albania     68.43292
 9   3      Algeria     59.03017
10   4       Angola     37.88350
11   5    Argentina     69.06042
12   6    Australia     74.66292
13   > lifeExp_bycountry %>% filter(meanlifeExp == max(meanlifeExp))
14   Source: local data frame [1 x 2]
15     country meanlifeExp
16      (fctr)        (dbl)
17   1 Iceland     76.51142
18   > lifeExp_bycountry %>% filter(meanlifeExp == min(meanlifeExp))
19   Source: local data frame [1 x 2]
20          country meanlifeExp
21           (fctr)        (dbl)
22   1 Sierra Leone    36.76917
```

**SLIDE** (Group by multiple variables)

- Talk around slide

- **Use multiple variables with** `group_by()` , `summarize()`

    - Can write this in the script for sanity

```
 1  > gdp_bycontinent_byyear <- gapminder %>%
 2  + group_by(continent, year) %>%
 3  + summarize(mean_gdpPercap=mean(gdpPercap))
 4  > head(gdp_bycontinent_byyear)
 5  Source: local data frame [6 x 3]
 6  Groups: continent [1]
 7    continent  year mean_gdpPercap
 8       (fctr) (int)          (dbl)
 9  1    Africa  1952       1252.572
10  2    Africa  1957       1385.236
11  3    Africa  1962       1598.079
12  4    Africa  1967       2050.364
13  5    Africa  1972       2339.616
14  6    Africa  1977       2585.939
15  > gdp_pop_bycontinents_byyear <- gapminder %>%
16    group_by(continent,year) %>%
17    summarize(mean_gdpPercap=mean(gdpPercap),
18             sd_gdpPercap=sd(gdpPercap),
19             mean_pop=mean(pop),
20             sd_pop=sd(pop))
21  > head(gdp_pop_bycontinents_byyear)
22  Source: local data frame [6 x 6]
23  Groups: continent [1]
24    continent  year mean_gdpPercap sd_gdpPercap mean_pop    sd_pop
25       (fctr) (int)          (dbl)        (dbl)    (dbl)     (dbl)
26  1    Africa  1952       1252.572     982.9521  4570010   6317450
27  2    Africa  1957       1385.236    1134.5089  5093033   7076042
28  3    Africa  1962       1598.079    1461.8392  5702247   7957545
29  4    Africa  1967       2050.364    2847.7176  6447875   8985505
30  5    Africa  1972       2339.616    3286.8539  7305376  10130833
31  6    Africa  1977       2585.939    4142.3987  8328097  11585184
```

**SLIDE** ( `mutate()` )

- Talk around slide

- `mutate()` **lets us create new variabls on the fly**

  - We can calculate total GDP from GDP per person, and population

```
 1  > head(gapminder %>% mutate(gdp_billion=gdpPercap*pop/10^9))
 2        country year       pop continent lifeExp gdpPercap gdp_billion
 3  1 Afghanistan 1952  8425333      Asia  28.801  779.4453    6.567086
 4  2 Afghanistan 1957  9240934      Asia  30.332  820.8530    7.585449
 5  3 Afghanistan 1962 10267083      Asia  31.997  853.1007    8.758856
 6  4 Afghanistan 1967 11537966      Asia  34.020  836.1971    9.648014
 7  5 Afghanistan 1972 13079460      Asia  36.088  739.9811    9.678553
 8  6 Afghanistan 1977 14880372      Asia  38.438  786.1134   11.697659
 9  > gdp_pop_bycontinents_byyear <- gapminder %>%
```

```
10  +     mutate(gdp_billion=gdpPercap*pop/10^9) %>%
11  +     group_by(continent,year) %>%
12  +     summarize(mean_gdpPercap=mean(gdpPercap),
13  +               sd_gdpPercap=sd(gdpPercap),
14  +               mean_pop=mean(pop),
15  +               sd_pop=sd(pop),
16  +               mean_gdp_billion=mean(gdp_billion),
17  +               sd_gdp_billion=sd(gdp_billion))
18  > head(gdp_pop_bycontinents_byyear)
19  Source: local data frame [6 x 8]
20  Groups: continent [1]
21    continent  year mean_gdpPercap sd_gdpPercap mean_pop    sd_pop mean_gdp_billion
22       (fctr) (int)          (dbl)        (dbl)    (dbl)     (dbl)           (dbl)
23  1    Africa  1952       1252.572     982.9521  4570010   6317450         5.992295
24  2    Africa  1957       1385.236    1134.5089  5093033   7076042         7.359189
25  3    Africa  1962       1598.079    1461.8392  5702247   7957545         8.784877
26  4    Africa  1967       2050.364    2847.7176  6447875   8985505        11.443994
27  5    Africa  1972       2339.616    3286.8539  7305376  10130833        15.072242
28  6    Africa  1977       2585.939    4142.3987  8328097  11585184        18.694899
29  Variables not shown: sd_gdp_billion (dbl)
30  > gdp_pop_bycontinents_byyear <- gapminder %>%
31  +     group_by(continent,year) %>%
32  +     summarize(mean_gdpPercap=mean(gdpPercap),
33  +               sd_gdpPercap=sd(gdpPercap),
34  +               mean_pop=mean(pop),
35  +               sd_pop=sd(pop)) %>%
36  +     mutate(gdp_billion=mean_gdpPercap*mean_pop/10^9)
37  > gdp_pop_bycontinents_byyear <- gapminder %>%
38  +     group_by(continent,year) %>%
39  +     summarize(mean_gdpPercap=mean(gdpPercap),
40  +               sd_gdpPercap=sd(gdpPercap),
41  +               mean_pop=mean(pop),
42  +               sd_pop=sd(pop)) %>%
43  +     mutate(mean_gdp_billion=mean_gdpPercap*mean_pop/10^9)
44  > head(gdp_pop_bycontinents_byyear)
45  Source: local data frame [6 x 7]
46  Groups: continent [1]
47    continent  year mean_gdpPercap sd_gdpPercap mean_pop    sd_pop mean_gdp_billion
48       (fctr) (int)          (dbl)        (dbl)    (dbl)     (dbl)           (dbl)
49  1    Africa  1952       1252.572     982.9521  4570010   6317450         5.724268
50  2    Africa  1957       1385.236    1134.5089  5093033   7076042         7.055054
51  3    Africa  1962       1598.079    1461.8392  5702247   7957545         9.112641
52  4    Africa  1967       2050.364    2847.7176  6447875   8985505        13.220489
53  5    Africa  1972       2339.616    3286.8539  7305376  10130833        17.091772
54  6    Africa  1977       2585.939    4142.3987  8328097  11585184        21.535946
```

# Creating publication-quality graphics

**SLIDE** (Creating publication-quality graphics)

**SLIDE** (Learning objectives)

- Talk around slide

## The grammar of graphics

**SLIDE** (The grammar of graphics)

- Talk around slide

- **Grammar of graphics is non-intuitive, but gives advantages**

    - Data and its representation handled separately
    - Means that components can be customised to a particular representation easily

**SLIDE** (A basic scatterplot)

- Talk around slide

```
1  > library(ggplot2)
2  > qplot(lifeExp, gdpPercap, data=gapminder, colour=continent)
```

- **Show the plot**

    - Describe features
    - x-, y-axes; colours by continent; legend
    - main features - Europe high life expectancy, Africa low GDP per capita

- **What is happening under the surface? How can you reproduce this?**

    - Convenience functions can be quick and easy, but aren't readily modifiable
    - We'd like to build plots *like* this in other situations - how can we do that?

**SLIDE** (What is a scatterplot? Aesthetics…)

- Talk around slide

**SLIDE** (What is a scatterplot? Aesthetics…)

- Talk around slide

- **Aesthetics decide where and how data are plotted**

    - They essentially create a new dataset that contains aesthetic information

**SLIDE** (What is a scatterplot? `geom` s)

- Talk around slide

- **`geom`s determine the "type" of plot**

    - Not all `geom`s make sense for a given dataset (though they may be 'grammatical')
    - Can combine multiple `geom`s to produce new graphs

**SLIDE** (`ggplot2` layers)

- Talk around slide

**SLIDE** (Building a scatterplot)

- **Creating a `ggplot` object**

    - Can't plot these directly
    - Can store them in variables for convenience/reproducibility

```
 1   > ggplot(data=gapminder, aes(x=lifeExp, y=gdpPercap))
 2   > p <- ggplot(data=gapminder, aes(x=lifeExp, y=gdpPercap))
 3   > > str(p)
 4   List of 9
 5    $ data        :'data.frame':    1704 obs. of  6 variables:
 6     ..$ country  : Factor w/ 142 levels "Afghanistan",..: 1 1 1 1 1 1 1 1 1 1 ...
 7     ..$ year     : int [1:1704] 1952 1957 1962 1967 1972 1977 1982 1987 1992 1997 ...
 8     ..$ pop      : num [1:1704] 8425333 9240934 10267083 11537966 13079460 ...
 9     ..$ continent: Factor w/ 5 levels "Africa","Americas",..: 3 3 3 3 3 3 3 3 3 3 ...
10     ..$ lifeExp  : num [1:1704] 28.8 30.3 32 34 36.1 ...
11     ..$ gdpPercap: num [1:1704] 779 821 853 836 740 ...
12     […]
```

- **We need to add a layer**

    - At minimum, use a `geom`
    - This uses the default dataset we specified in `p`, unless told otherwise
    - `geom_point` tells `ggplot2` we want to represent data as points (scatterplot)
    - We get only a scatterplot of points, but no colours

```
 1   > p + geom_point()
```

- **We can modify aesthetics**

    - In the default dataset, or in the `geom` layer
    - Aesthetics/data in the `geom` layer override those in the default

```
 1   > p + geom_point(aes(colour=continent))
 2   > p <- ggplot(data=gapminder, aes(x=lifeExp, y=gdpPercap, colour=continent))
 3   > p + geom_point()
```

**SLIDE** (Challenge 1)

Solution:

```
1   > p <- ggplot(data=gapminder, aes(x=year, y=lifeExp, colour=continent))
2   > p + geom_point()
```

This is not a good way to view the data - we need a new geometry!

# Layers

**SLIDE** (Layers)

- Talk around slide

- **The last challenge representation didn't look good**

    - Change `geom` to line chart

```
1   > p + geom_line()
```

- This looks wrong

    - Lines connect continents, not countries (which is what we want)

- **Group data on a variable**

    - Use `by` to group data by country

```
1   > p + geom_line(aes(by=country))
```

- That looks better

- **Overlay a second `geom` to see datapoints**

    - Use the `+` operator to keep adding `geom`s
    - Layers are drawn in the specified order

```
1   > p + geom_line(aes(by=country)) + geom_point()
2   > p + geom_line(aes(by=country)) + geom_point(aes(colour=NULL))
3   > p + geom_point(aes(colour=NULL)) + geom_line(aes(by=country))
```

# Transformations and statistics

**SLIDE** (Transformations)

- Talk around slide

- **Scaling axes**

  - Difficult to distinguish GDP on the y-axis
  - Rescale with a transformation

```
1  > p <- ggplot(data=gapminder, aes(x=lifeExp, y=gdpPercap, colour=continent))
2  > p + geom_point()
3  > p + geom_point() + scale_y_log10()
```

- **Transformations can be layered**

```
1  > p + geom_point(aes(size=pop)) + scale_size("population")
2  > p + geom_point(aes(size=pop)) + scale_size("population") + scale_y_log10()
```

- **Scaling colours**

  - Transformations are also how colours are 'scaled'

```
1  > p + geom_point() + scale_y_log10() + scale_colour_brewer()
2  > p + geom_point() + scale_y_log10() + scale_colour_grey()
```

**SLIDE** (Statistics)

- Talk around slide

- **Adding a smoother to the data**

  - Adds as another layer on the plot

```
1  > p <- ggplot(data=gapminder, aes(x=lifeExp, y=gdpPercap))
2  > p + geom_point()
3  > p + geom_point() + scale_y_log10()
4  > p + geom_point() + scale_y_log10() + geom_smooth()
```

- **Adding a KDE**

  - Adds as another layer on the plot

```
1  > p + geom_point() + scale_y_log10() + geom_density_2d()
```

# Multi-panel figures

**SLIDE** (Multi-panel figures)

- Talk around slide

- **Faceting**

  - Grouping data by country, colouring by continent

- One big plot is messy, hard to read.
- Using `facet_wrap` splits out plots on groups

```
1  > p <- ggplot(data=gapminder, aes(x=year, y=lifeExp, colour=continent, by=country))
2  > p + geom_line()
3  > p + geom_line() + scale_y_log10()
4  > p + geom_line() + scale_y_log10() + facet_wrap(~continent)
```

- **Grouping on country**

  - Even the continent plots are a bit jumbled
  - Group by country just by changing the argument

```
1  > p + geom_line() + scale_y_log10() + facet_wrap(~country)
```

- Very hard to read in `RStudio`
- Export graph as pdf and visualise
  - Click `Export -> Save as PDF`
  - PDF Size: A4
  - Orientation: Landscape
  - File name (something sensible)
  - View plot after saving
  - `Save`

**SLIDE** (Challenge 2)

Solution:

```
1  > p <- ggplot(data = gapminder, aes(x = gdpPercap, fill=continent))
2  > p + geom_density()
3  > p + geom_density(alpha=0.6)
4  > p + geom_density(alpha=0.6) + scale_x_log10()
5  > p + geom_density(alpha=0.6) + scale_x_log10() + facet_wrap(~year)
```

# Wrapping up

**SLIDE** (Wrapping Up)

**SLIDE** (Learning objectives)

- Talk around slide

**SLIDE** (Best practices)

- Talk around slide