

# 1 R Programming

## List

- `[[idx]]`: get element in a list
- `str(ls)`: get **structure** of a list (similar to summary)
- `saveRDS` and `loadRDS`
- `unlist`: convert list to vector **[IMPT]**

## Recycling Rule

- shorter vectors are recycled until they match the length of the longest vector
- the length of the longest vector must be a multiple of the shorter vector in arithmetic operations!

## Useful functions

- `sample(x, size, replace, prob)`
  - `size`: length of output vector
  - `replace`: if TRUE, then sampling is with replacement
  - `prob`: a vector of probability weights
- `any(duplicated(vec))`: returns true or false if there are any duplicated elements in a vector
- `rep(x, times, length.out)`
- `table()`
- `args(func)`: list the arguments of a function
- `seq(from, to, by, length)`
- `paste(v1, v2, sep)`: concatenate vectors after converting them to characters
  - `sep`: separator between elements of `v1` and `v2`
  - The recycling rule applies when `length(v1) != length(v2)`
- apply function family: apply function to each row (1) or column (2)
  - `apply(X, margin, func, ...)`
    - \* Note that `X` must be a **matrix** or **df** in `apply`
  - `sapply` returns a vector or a matrix, **input must be 1 dimensional!**
  - `lapply` returns a list, useful when the output of the function may not be all of the same length/-type, **input must be 1 dimensional!**
  - `replicate(n, func)`: replicate anonymous function `n` number of times (especially useful for random number generations)
  - `tidy::tidyselect::group_by()`: used to apply function and then group them into a table using grouping index
  - `mapply(func, arg1, arg2, arg3, ...)`: like `sapply` but takes multiple vectors containing arguments to `func`

## Function debugging

- `cat("...")`: used to print statements
- `browser()`: debugging with breakpoint

## Important classes

### Strings

- Start by importing `tidyverse` and `stringr`

- Library functions
  - `str_length`: returns vector of string lengths
  - `str_c(..., sep)`: concatenate strings with optional separator
  - `str_sub(string, start, end)`: returns vector of substrings
- Regular expressions (`str_view()` to test out regex), [Tidyverse Article](#)
  - to match an **a** at the beginning of a string  
`str_view(x, "^a")`
  - to match an **a** at the end of a string  
`str_view(x, "a$")`
  - to match an **a** or **e** at the end of a string  
`str_view(x, "[ae]$")`
  - to match a string of 3 chars with **a** in the middle  
`str_view(x, ".a.")`
- `str_detect(vec, regex)`: returns a boolean vector
  - `|`: means or  
`str_detect(street_names, "Jurong|Boon Lay")`
  - `+`: means modifier (pattern detected 1 or more times)
  - `()`: to group stuff
  - `\\w`: any word
  - `[0-9]`: can be 0 to 9
  - `\\d`: any number
    - \* `\\d{3,6}` to search for digits repeating between 3 and 6 times
  - **[IMPT]** `?about_search_regex` for help
  - **[IMPT]** `?base::regex :help` for regex from R base package; `[:punct:]`, `[:digit:]`, `[:space:]`
- `str_extract(vec, regex)`: returns a vector of strings, particularly helpful for `".a."` regex

```
1 # To find the number of eggs given
  a sentence
2 str_extract(sent, "[0-9]+(?= eggs)")
3 # ?= is a look behind operator
4 # ?<= is a look ahead operator
5
```

- `str_trim`: to trim trailing whitespaces
- `str_split`
- `str_replace`

```
1 # to remove duplicate words
2 str_replace(sent_type, "\\b(\\w+)\\b \\1", "\\1")
```

Note that `\\b` means word boundary and `\\1` means group boundary 1

- `str_match`

**[IMPT]** USE `vignette('stringr')` and `vignette('regular-expressions')` for help

- `devtools::install_github("gadenbuie/regexpal")` to install `regexpal` GUI, need to install `devtools` library first
- Also Tools → Addins → Browse Addins.. → `regexpal` (cheatsheet/GUI)

## Factors

`factor(vec, levels=c(...))`: convert `vec` to factors with fixed levels

`unique(vec)`: returns a vector with unique values

## Date

- **[IMPT]** `?strptime` for help page
- `as.Date(x, format)`: convert string `x` to Date object  
e.g. `as.Date("2014/02/22", "%Y/%m/%d")`
- `months(d)`: what month of the year is the date in?
- `weekdays(d)`: what day of the week is the date on?
- `Sys.Date()`
- `cut(x, breaks, labels)`: usually used to group dates that fall into a month/week/quarter
  - `breaks`: numeric vector/string ("month", "week")
  - `labels`: if TRUE, return a label vector
- `seq(d, d+365, by="1 week" or "1 quarter")`

## Basic Plotting

### plot()

- `pch`: abbr. for plotting character
- ```
1 # show all pch characters
2 example(pch)
```
- `col`:
- ```
1 # show all preset colours
2 colours()
3 # set custom colour, alpha is transparency
4 col <- rgb(..., alpha=?)
```
- `cex`: abbr. for character expansion
  - `bty`: change box borders
  - **[IMPT]** `?par` shows all parameters for `plot()`
  - use `points()` or `lines()` to add more stuff to an existing plot
    - `segments(x_)`

### barplot()

- `horiz=TRUE` flip y and x axes
- `las` (under `?par`)
- `par()`: **[IMPT]** lists all the default parameters for plots (`mar`, `mfrow` etc.)
- How to set graphical param?

```
1 # 1 row 2 columns plot
2 opar <- par(mfrow=c(1,2))
3 # plot some stuff
4 par(opar) # to set it back to default
5
```

### hist()

- `freq`: makes the y-axis a proportion of all the total shit (count/total), not total count using integer

## 2 Stringr

### (to convert numeric to string) Fixed vs scientific format

- Scientific: 1.989e+30 to denote  $10^{30}$
- `format(x, scientific=TRUE)` to format number to string by specifying digit numbers etc.
- **[IMPT]** `digits=` will format the smallest number so that it only has the specified significant digit, and other numbers in the vector follows

```
1 format(c(0.0011, 0.011, 1),
2        digits=1)
3 > [1] "0.001" "0.011" "1.000"
```

- `formatC(x, format="f" OR "e" or "g")`  
f stands for fixed, e for scientific, and g for scientific if it saves space

### Stringr functions

- `str_c`: concatenate like paste
- `str_length`: find length
- `str_sub`
- `str_detect`: returns boolean vectors
- `str_subset`:
- `str_count`
- `str_split`: `n=` returns maximum number of n elements, `simplify=` returns a matrix
  - **[IMPT]** `type=boundary("sentence")`
- `str_match`: returns a matrix with the capture or () regex
- `str_to_upper()`: returns a vec with all uppercase elements
- `str_to_lower()`
- `regex(expr, ignore_case = TRUE)`: tells regex to ignore case

### Rebus package

- `install.packages("rebus")` ⇒ `library(rebus)`
- rebus syntax can be used for stringr pattern instead of regex

```
1 pattern = START %R% "a"
2 # strings that start with "a"
3 # same as regex "^a"
4 # END is also possible
5 # %R% is read as 'then'
```

- `ANY_CHAR`
- `WRD`: word, `SPC`: Space

```
1 # to capture word ending in ING
2 one_or_more(WRD) %R% "ING"
3 # equals to \w+ING
4
```

- `or(p1, p2)`: kinda like | in regex
  - `or1(vec)`: pass vec as alternatives instead of arguments
- `char_class("Aa")`: kinda like "[Aa]" in regex
- `negated_char_class("aieuoAIUEO")`: self-explanatory
- `optional()`: ? in regex
- `zero_or_more()`: \* in regex
- `one_or_more()`: + in regex
- `repeated()`: {m,n} in regex

- `exactly()`: matches exact string
- `capture(pattern)`: group parts of pattern together, which is `()` in regex format
  - \*use REF1, REF2, REF3 to refer to the capture group (exact match) which is `\1`, `\2` and so on in regular regex

### String functions

- `stri_isempty()`: returns boolean

### Miscellaneous

- `strftime(date, format)`: string from time object
- `as.POSIXct(date_string, format)`: convert string to Date time

## 3 R Markdown (RMD)

- .yaml header

```

1  title: "...
2  output:
3    html-document:
4      toc: true #table of content
5      toc_float: true # floating TOC
6      at the left side of the window
7      collapsed: true
8      smooth_scroll: true
9      toc_depth: 2
10     number_sections: true/false
11 date: 'r format(Sys.time(), "%d %
12       B %Y")'
13 params:
14   country: Indonesia

```

- how to reference?? ⇒ I want die liao `'r params$country'`
- Referencing is important as it allows more control over the report, don't need to manually change the name of every variable if we want something else

- R Setup **[IMPT]**, will apply settings globally

```

1  '{r setup, include=FALSE}
2  knitr::opts_chunk$set(fig.align='
3  center', echo=TRUE)

```

- Use `'r var'` to insert inline code and ask R to run it
- Figure
  - `include=FALSE/TRUE`: to include the output or not
  - `fig.width, fig.height, fig.dim = c(w,h), out.width="XX%"`
  - `fig.align='left'/'centre'`
  - `fig.cap` for captions
- Bulleted list: just indent and use `'-'`
- Display table: use `kable(df, col.names=c(...))`
  - Important parameters: `caption`, `align="ccc"` or `"lll"` for text alignment inside boxes

### Code Chunk Settings

- `include=FALSE` doesn't print the code

- `echo=FALSE` usually for plots, don't include the actual code but just runs it
- `eval=FALSE` code chunk is not run/evaluated
- `collapse=TRUE` combines text output and source code in single block
- `message=FALSE`
- `warning=FALSE`
- `error=TRUE` will continue to knit the file even when there are errors and will include error messages in the file

## 4 Importing Data

**[IMPT]** use `read.delim` or `readLines` if none is working

### CSV Files

`read.csv()`: main arguments:

- `file`: filename/path
- `skip`: skip lines?
- `header`: default is TRUE
- `row.names`
- `stringsAsFactors`
- `na.strings`: what are the NA values
- `colClasses`: what classes are the columns (in terms of class names vector)

### Procedure when dealing with CSV:

- `apply(salaries, 2, function(x) sum(is.na(x)))`
- **[IMPT]** (check if any column has missing values)
- if `read.csv` doesn't work, can try `readLines` and `str_split` to split commas

### Excel Files

- import `readxl`, data is in the form of a tibble
- `read_excel(path, sheet=?)`: `sheet` parameter can be string or integer
- `sheet_names(path)`: to retrieve sheet names

### JSON Files

- import `jsonlite`
- `fromJSON(txt)`: takes up text/string object as an argument
- `readLines(path)`: returns a string **[IMPT]** line break will count as another element of a vector
- `prettyfy()`
- **[IMPT]** How to convert list to data frame?
  1. create a function `ls_to_df` which returns `data.frame` given an element of a list
  2. lapply the list to return a list of dataframes
  3. use `do.call` to combine the individual dataframes into one single dataframe

```

1 df_row_list <- lapply(list, ls_
2   to_df)
3 # combine repeatedly
4 do.call(rbind, df_row_list)

```

- Some thoughts **[IMPT]** Are there missing data for any observation?? if yes then remove

## 4.1 OOP in R

**[IMPT]** Main purpose: call function the same way (with similar syntax but different behaviour for each class) e.g. plot works differently for timeseries and vectors **S3 classes**

- methods: to search for available methods
- summary

```
1 studentBio <- list(studentName = "Harry
  Potter", studentAge = 19,
  studentContact="London")
2 class(studentBio) <- "StudentInfo"
3
4 # how to assign method
5 contact <- function(object) {
6   UseMethod("contact")
7 }
8 contact.StudentInfo <- function(object)
9 {
10  cat("Your contact is", object$
11    studentContact, "\n")
12 }
13 # can just call contact(studentBio)
14 without .StudentInfo
```

### S4 classes

```
1 # How to set Class with slots
2 setClass("employee", slots=list(name="
  character", id="numeric", contact="
  character"))
3
4 # Constructor
5 obj <- new("employee", name="Steven", id
6   =1002, contact="West Avenue")
```

**[IMPT]** How to add method?

```
1 setMethod("show",
2   signature(object="employee"),
3   definition=function(object) {
4     # do stuff
5   })
```

**[IMPT]** Tips for dealing with S4 data

- isS4(obj): check if obj is S4
- slotNames(obj) list all the attributes/slots
- methods(class="????"): to list out all the methods
- methods(generic.function="plot"): to list out all the classes a method can be applied to
- vignette("class"): for documentation

### RC classes

## 5 Databases

### How to connect?

- Install the requisite package on R
- Authenticate to the database server
- Query/Extract the data
- Analyse the data
- Close the connection

## 5.1 MongoDB

### Steps to connect

- **[IMPT]** MongoDB Tutorial Docs
- Code to connect

```
1 library(mongolite)
2 library(jsonlite)
3 # eXXXXX:pwd
4 credentials <- paste0(readLines("
  mongo_user_pwd.txt", warn=FALSE)
5   , collapse=":")
6 connection_string <- paste0("
  mongodb://", credentials, "
  @rshiny.nus.edu.sg:2717/test")
7 con2 <- mongo(verbose=TRUE,
8   collection="restaurants", url=
9   connection_string)
10 con2$count()
```

## 5.2 Data from Web

### 5.2.1 Download File from Link

- how to download

```
1 imda_url <- "https://data.gov.
2 sg/dataset/02c1f624-489f-40ad-8
3 fdd-5e66e46b2722/download"
4 return_val <- download.file(
5   imda_url, "../data/imda_data.zip"
6 )
7 con <- unz("../data/imda_data.
8 zip", "wage-02-size2-annual.csv"
9 )
10 wages_data <- read.csv(con,
11   header=TRUE)
```

- download.file(), mode="wb" for Windows
- unz: to unzip

### 5.2.2 Developer API

- Normal browser  $\xleftrightarrow[\text{response}]{\text{request}}$  Web server
- request data from server that is continuously running
- **[IMPT]** Usually for Real-time data
- how to get data?

```
1 library(httr)
2 set_config(verbose())
3 url <- "https://api.data.gov.sg/v1/
4 transport/taxi-availability"
5 taxi_avail <- GET(url, query=list(
6   date_time="2022-08-01T09:00:00"
7 ))
8 taxi_data <- content(taxi_avail)
```