

# Introduction to R

Day 1

Transforming Publishing

[phs.transformingpublishing@phs.scot](mailto:phs.transformingpublishing@phs.scot)

# Pathway

## Intro

Intro to Data and Tools, Overview of R

## Foundations

Commenting, Types, Variables, Statements, Data Structures, Packages

## RStudio

Desktop/Server Version, Interface, Customisation, R Scripts, Hints/Tips

## Workflow

Overview (Collect, Explore, Wrangle, Viz, Outputs), Git(Hub/ea), RMarkdown, RAP, Templates, Style Guide

## Wrangle

Tidyverse (dplyr/magrittr), Pipes, Functions (Filter, Mutate, Arrange, etc.), PHS Methods

## Explore

Mean, Median, Summary Function, Frequencies/Cross-Tabs

## Data Flow

Directories/File Paths, CSVs, SPSS (haven), SMRA/Other Databases

## Visualise

Intro to ggplot2, Line Graphs, Bar Plots, Scatterplots, Customisation

## Output

Overview of RMarkdown, Shiny, etc.

## Review

Overview, Next Steps, Q&A



**Introduction – the why**

# R

- is a **programming language** widely used for *data analysis*, statistics, and *graphics*;
- is **open source**, available on all major operating systems;
- has the functionality to go from raw data to interactive reports and web apps;
- and it's part of the **PHS analytical strategy**.



**Foundations – the what**

# Building Blocks

- **Basic Data Types** – how is fundamental data, like numbers and text, represented?  
This is then the foundations of more complex, composite data types, e.g. tables.
- **Variables** – named storage to track "objects" across our program
- **Statements** – a complete line of code, made of expressions and operators
- **Control Flow** – branching (if statements) and iteration (loops)
- **Functions** – reusable code that *can* take inputs and give outputs

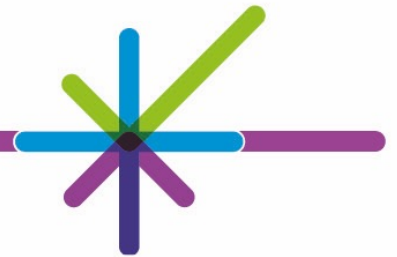


# Basic Data Types

- Character (String) – e.g. "Hello World!"
- Numeric (Float/Real) – e.g. 123.5
- Logical (Boolean) – e.g. TRUE

```
typeof("Hello World")  
is.numeric(123.5)  
print(typeof(TRUE))
```

```
> [1] "character"  
> [1] TRUE  
> [1] "logical"
```



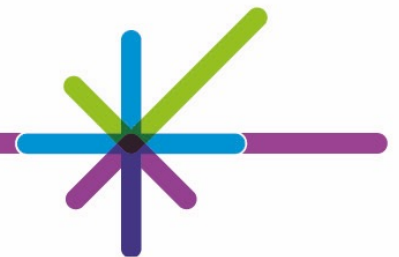
# Type Conversion

`as.<data_type>()`

- `as.character()` conversions tend to succeed without fault
- `as.numeric()` – TRUE and FALSE become 1 and 0, character types need to be formatted correctly
- `as.logical()` – everything except 0 becomes TRUE for numeric conversions, character can be upper, lower, or proper case versions

```
as.character(123.5)
as.numeric("123.5")
as.logical("False")
```

```
> [1] "123.5"
> [1] 123.5
> [1] FALSE
```





# Variables

- **Naming** – letters, numbers, dots `.` or underscores `\_` are all okay. However, you **can't** start with an underscore or number, or a dot then a number. Any existing terminology is also reserved from being used as a variable.  
*Following style guidance is also important.*
- **Assignment** – variables are assigned mainly with `<-` but you may also see `=` being used.

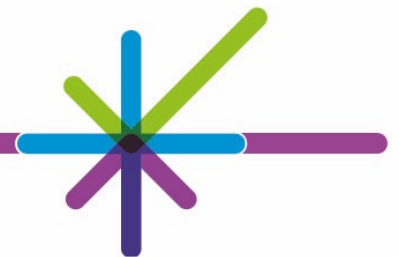
```
# Good
Totals
sumOfPatientsUnder60

# Bad
60YearOldPatients
_template
TRUE
```



# Operators

Precedence	Operator	Description
1	^	Exponentiation (right to left evaluation)
2	%	Modulus
3	*    /	Multiplication, Division
4	+    -	Addition, Subtraction
5	<   >   <=   >= ==   !=	Comparison Operators (Less Than, More Than, Less Than or Equal To, More Than or Equal To, Equal To, Not Equal To)
6	!	Logical NOT
7	&   &&	Logical AND
8		Logical OR



# Knowledge Check

We're going to use an app for a lot of our interactive work, especially today. (Hint: some questions have hint buttons).

[scotland.shinyapps.io/phs-rtraining-intro/](https://scotland.shinyapps.io/phs-rtraining-intro/)

- Foundations

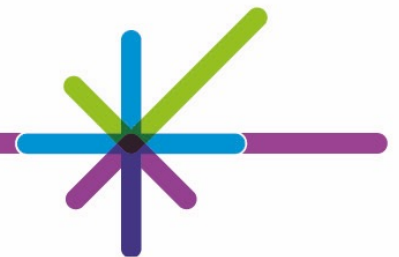


# Anatomy of a Program

```
# Example 1  
hello_world <- "Hello World"  
print(hello_world)
```

```
> [1] "Hello World"
```

- `# Example 1` – comment
- `hello_world` – variable
- `<-` – assignment operator (alt + -)
- `"Hello World"` – character (" or `)
- `print()` – function



# Style Guide

- **Naming** – variables and filenames should have meaningful names in snake\_case format, preferring all lower case.
- **Structure**
  - Space after a comma
  - No spaces before or after parenthesis
  - Comments to explain code and create sections within the code
  - Prefer `"` over `'` for characters

```
# Bad
pts <-c ( 'Al','Bert' )

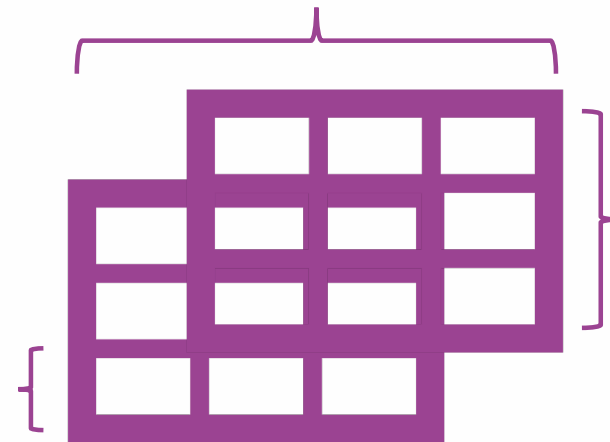
# Good
patients <- c("Al", "Bert")
```



**Foundations – the what**

# Data Structures

- Vectors
  - Lists
  - Factors
  - Matrices
  - Data Frame
- 
- `str()` provides an overview and description of the data structure



# Vectors

contain multiple objects of the *same* basic class

- Create: `c(...)` or `vector(<type>, <length>)`
- Access: `<vector>[<index>]`

```
c("a", "c", "f", "b")[1]  
c(2, 5, 1, "abc")[3:4]  
vector("logical", 4)
```

```
> [1] "a"  
> [1] "1" "abc"  
> [1] FALSE FALSE FALSE FALSE
```





# Knowledge Check

[scotland.shinyapps.io/phs-rtraining-intro/](https://scotland.shinyapps.io/phs-rtraining-intro/)

- Data Structures - Vectors



# Lists

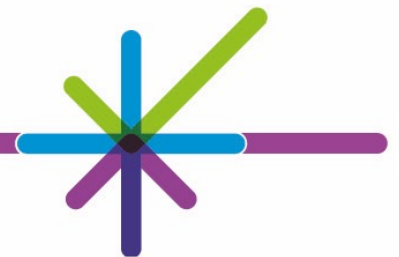
are a special type of vector that can contain objects of *different* classes, including other lists

- Create: `list(...)`
- Sub-list: `<list>[<index>]`
- Access: `<list>[[<index>]]`

```
list("abc", 4, FALSE)[1:2]  
list(list(2, 3), "abc")[[2]]
```

```
> [[1]]  
> [1] "abc"  
> [[2]]  
> [1] 4
```

```
> [1] "abc"
```



# Naming Elements

can be done on vectors and lists during or after creation.

- At creation:

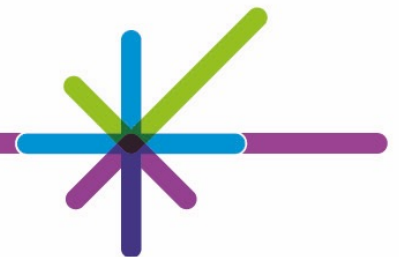
```
c("<name>", = <item>) or  
list("<name>", = <item>)
```

- After:

```
names(<object>) <-  
c("<name>")
```

```
x <- list("Ch" = "a", "Nm" = 2)  
names(x) <- c("Char", "Num")  
x$Char
```

```
> [1] "a"
```



# Factors

are used to represent categorical data, with both ordered and unordered variations.

- **Create:** `factor(c(...), ordered = <TRUE/FALSE>)`

- **Levels:** `factor(c(...), levels = c(...))`

```
x <- factor(c("M", "F", "M"),  
            levels = c("F", "M"))  
x
```

```
> [1] M F M  
Levels: F M
```



# Knowledge Check

[scotland.shinyapps.io/phs-rtraining-intro/](https://scotland.shinyapps.io/phs-rtraining-intro/)

- Data Structures - Factors



# Matrices

expand our dimensions with a `nrow` and `ncol` arguments, constructed column-wise

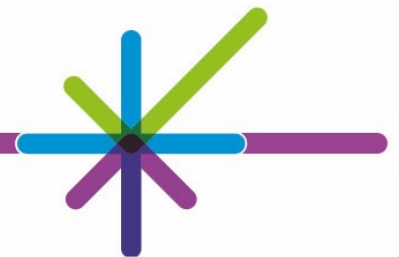
- **Create:** `matrix(<data>, nrow = <int>, ncol = <int>)`

- **Access:**  
`<matrix>[<row>, <col>]`

```
x <- matrix(1:6, 2, 3)
x
x[2, 3]
```

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

```
> [1] 6
```



# Data Frames

are used to store tabular data, each column contains one variable, each row contains an observation

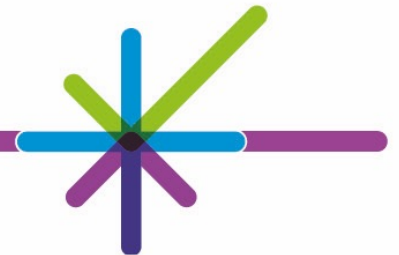
- Create: `data.frame("<name>" = <element(s)>)`

- Subset: `[]`

- Access: `[[ ]]` or `$`

```
data.frame(name = c("Harry",  
                    "Sarah"),  
           score = c(62, 91))
```

```
>      name  score  
> 1  Harry    62  
> 2  Sarah    91
```



# Knowledge Check

[scotland.shinyapps.io/phs-rtraining-intro/](https://scotland.shinyapps.io/phs-rtraining-intro/)

- Data Structures - Data Frames





# Tibbles

data frames that attempt to make our lives a bit easier. First, load the package:

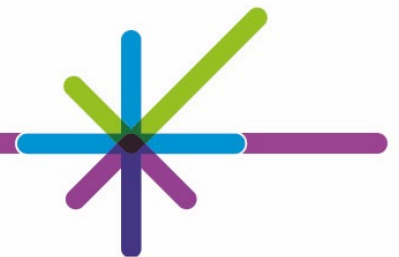
```
library(tibble)
```

- **Create:** `tibble("<name>" = <element(s)>)` or coerce an existing data frame with `as_tibble(<object>)`

- **Subset and Access:** `[]` `[[ ]]` or `$`

```
tibble(name = c("Harry",  
                "Sarah"),  
       score = c(62, 91))
```

```
>      name  score  
> 1 Harry    62  
> 2 Sarah    91
```



**Foundations – the what**

# Anatomy of a Function

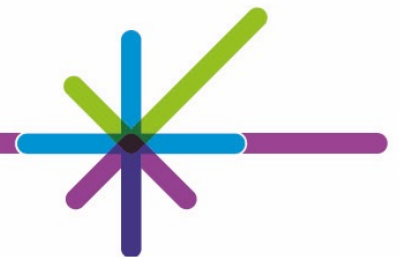
Functions allow us to bundle code for reuse, taking inputs, doing something and, optionally, providing outputs.

```
<name> <- function(<inputs...>) {  
  <code>  
  return(<outputs...>)  
}
```

```
mult_2 <- function(x) {  
  x <- x * 2  
  return(x)  
}
```

```
mult_2(4)
```

```
> [1] 8
```



# Packages

are used to expand the functionality of R with more functions

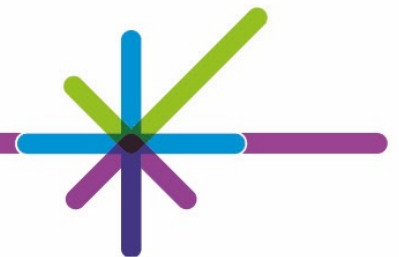
- **Install:**

```
install.packages("<package>")
```

- **Load:** `library(<package>)` or `require(<package>)` if loading packages as part of functions as it returns a logical value and a warning if the package isn't installed.

```
install.packages("tidyverse")  
library(tidyverse)
```

*- output varies by package, warnings (not errors) are normal. An example would be where a function 'masks' that of another R function.*



## Control Flow - if statements

Package: dplyr

Load package, `library(dplyr)`

- `if_else(<condition>, <true>, <false>)`

```
library(dplyr)
x <- 5
if_else(x > 10, TRUE, FALSE)
```

```
> [1] FALSE
```



# Knowledge Check

[scotland.shinyapps.io/phs-rtraining-intro/](https://scotland.shinyapps.io/phs-rtraining-intro/)

- Control Flow - If



## Control Flow - case statements

Package: dplyr

Load package, `library(dplyr)`

- `case_when(<condition> ~ <result>)`

```
library(dplyr)
x <- c(1, 2, 3, 4, 5)
case_when(x < 3 ~ "LT3",
          x %% 2 == 0 ~ "Even")
```

```
> [1] "LT3" "LT3" NA "Even" NA
```



# Knowledge Check

[scotland.shinyapps.io/phs-rtraining-intro/](https://scotland.shinyapps.io/phs-rtraining-intro/)

- Control Flow - Case





## Iteration – for loop

allowing us to do the same thing repeatedly with different inputs.

- `for(<value> in <sequence>)`  
  `{<statement>}`

```
files <- list.files(pattern = ".csv")
all_files <- list()
for(i in seq_along(files)) {
  all_files[[i]] <-
    read.csv(files[i])
}
```



## Iteration – loop with purrr

Package: purrr


Load: `library(purrr)`

- `map(<object>, <function>)`






```
library(purrr)
files <- list.files(pattern = ".csv")
all_files <- map(files, read_csv)
```




**RStudio**



FileEditCodeViewPlotsSessionBuildDebugProfileToolsHelp



Go to file/function


 Addins

Project: (None)R 3.5.1

Console~/

R version 3.5.1 (2018-07-02) -- "Feather Spray"  
Copyright (C) 2018 The R Foundation for Statistical Computing  
Platform: x86\_64-pc-linux-gnu (64-bit)  
  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.  
  
Natural language support but running in an English locale  
  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
  
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.  
  
> |






EnvironmentHistoryConnections

 Import Dataset

Global Environment

<input type="checkbox"/>	Name	Type	Length	Size	Value
Environment is empty					

FilesPlotsPackagesHelpViewer

 New Folder Upload Delete Rename More

Home

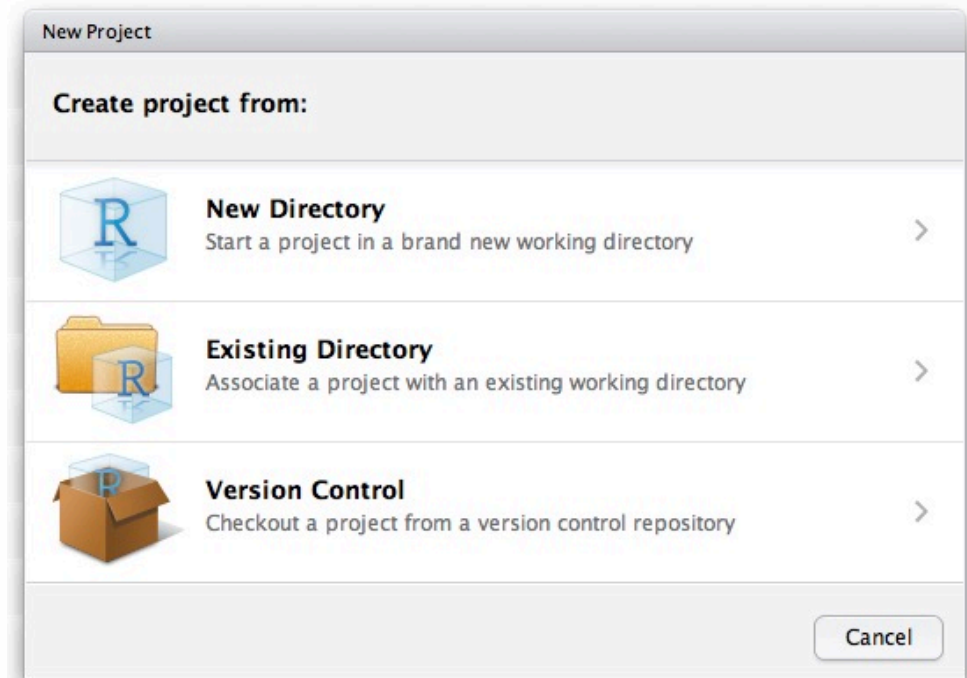
<input type="checkbox"/>	Name	Size	Modified
<input type="checkbox"/>	.Renvi	394 B	

# R Projects

keeps work separate, giving a project its own working directory, workspace, and history.

Opening an .Rproj file will:

- Start a new R session
  - Load project specifics and settings
  - The project directory is set as the current working directory.
- 
- **Create:** available in the Projects menu or global toolbar.



**Data Flow**

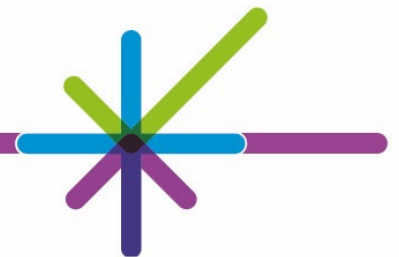
## Working Directory

- Current: `getwd()`
- Set new: `setwd(<filepath>)`
- We can also use the `here` package, with `here()`

RStudio also provides options through the user interface for navigating files and directories.

```
setwd("/home/learnr01/intro_R")  
getwd()  
here()
```

```
> [1] "/home/learnr01/intro_R"  
> [1] "/home/learnr01/intro_R"
```



# Read CSV

Package: readr

1. Load package, `library(readr)`
2. `read_csv(<filepath>)`
3. Check output

```
library(readr)
borders_csv <-
  read_csv("data/Borders.csv")
View(borders_csv)
```





# RDS

Package: readr

Load package, `library(readr)`

Read

- `read_rds(<filepath>)`

Write

- `write_rds(<object>, <filepath>)`

```
library(readr)

borders_RDS <-
  read_rds("data/borders.rds")

write_rds(borders,
  "data/borders.rds")
```



# Read SPSS

Package: haven

1. Load package, `library(haven)`
2. `read_sav(<filepath>)`
3. Check output

```
library(haven)
borders_spss <-
  read_sav("data/Borders.sav")
View(borders_spss)
```



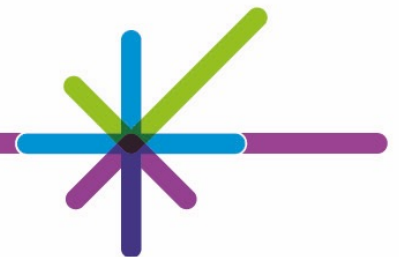
## Read Web

The packages/functions used will vary depending on the structure of the data. This example uses a CSV so the process to follow is the same as before.

Package: readr

1. Load package, `library(readr)`
2. `read_csv(<filepath>)`
3. Check output

```
library(readr)
hospital_codes <- read_csv("
  https://www.opendata.nhs.scot/dataset/[...].csv")
View(hospital_codes)
```



# Open Data – CKAN API

Package: ckanr

1. Load package, `library(ckanr)`
2. Set up connection and resource ID
3. `dplyr::tbl(<src>, <res>)`  
`%>%`  
`as_tibble()`
4. Check output

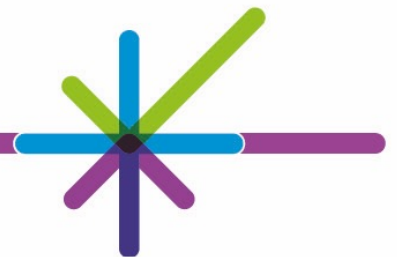
```
library(ckanr)
ckan <-
  src_ckan("https://www.opendata.nhs.scot")
res_id <- "<ID>"
resource <- dplyr::tbl(src =
  ckan$con, from = res_id) %>%
  as_tibble()
```



## Database (SMRA)

Package: odbc

1. Load package, `library(odbc)`
2. Connect: `smra_connection <- dbConnect(drv = odbc(),  
dsn = "SMRA",  
uid = .rs.askForPassword("SMRA Username:"),  
pwd = .rs.askForPassword("SMRA Password:"))`
3. Extract: `smr01 <- dbGetQuery(smra_connection,  
paste("<sql>"))`
4. Check output



# Write CSV

Package: readr

1. Load package, `library(readr)`
2. `write_csv(<object>, <filepath>)`
3. Check output

```
library(readr)
write_csv(borders,
          "data/borders.csv")
```

*- the write functions expect a dataframe as the object.*



**Explore**

## Mean/Median & Summary

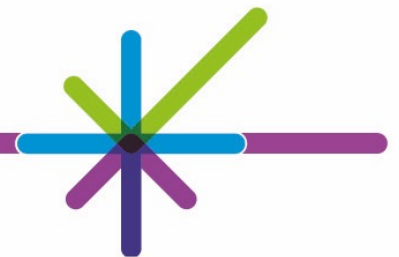
- `mean()` and `median()` are passed arrays of values (usually from a data frame) to return the mean and median value

```
mean(borders[["LengthOfStay"]])  
summary(borders$LengthOfStay)
```

- `summary()` returns all summary statistics based on a given array (usually from a data frame)

```
> [1] 4.297008  
>   Min. 1st Qu. Median   Mean 3rd Qu.  Max  
0.000  0.000   1.000  4.297  4.000   458
```

- We access columns using `[[]]` or `$`



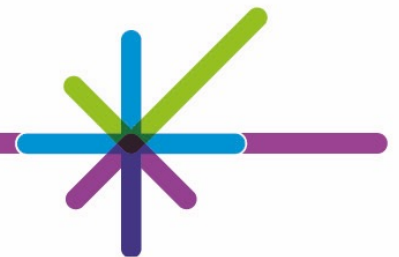


## Frequencies & Crosstabs

- Frequency: `table(<df>$<col>)`
- Crosstab: `table(<df>$<col1>, <df>$<col2>)`
- To add column and row totals, the function `addmargins()` can be used

```
addmargins(table(borders$HospitalCode, borders$Sex))
```

	1	2	3	Sum
A210H	1	0	0	1
B102H	56	100	0	156
B103H	50	108	0	158
...				
Sum	11947	13340	2	25289



## Exercise 1

1. Read in "Borders.csv" (giving the data frame an appropriate name)
2. What are the **mean, median, and max** values from the LengthOfStay variable? Can you do this in one step?
3. Produce a **frequency table** to check the sex variable, save this as an object with an appropriate name
4. Export the frequency table as a **csv** file. You'll need to use `as.data.frame()`



**Review**

# Project

1. Get familiar with RStudio and change some settings to your preference.
2. Create a new script and import some open data (using the API or .csv).
  - e.g. COVID-19 Daily Case Trends
3. Produce some summary statistics and a frequency table.



# Getting Help

- Vignettes (Help) / ``?<function>``
- Google / Stack Overflow  
tag queries "[r] & [tidyverse]"
- [R User Group Teams](#) – [Technical Queries](#)
- [Transforming Publishing](#)

