

Intro to R ENVX3003

Willem Vervoort

2019-06-12

Contents

1	Introduction	2
2	R as a modelling environment	3
2.1	R and R Studio	3
3	Working folder structure	4
3.1	The working directory	4
4	BASIC R	5
4.1	R as a calculator	5
4.2	Objects in R	5
5	STATISTICAL ANALYSIS AND DATA MANIPULATION	6
5.1	Packages to use	7
5.2	A dataframe in tidyverse	7
5.3	Reading data from different sources	10
5.4	More data manipulation (using tidyverse)	11
5.5	Different data set	13
6	PLOTTING	16
6.1	Exercise	20
7	Programming: if else and for loops	21
7.1	The “for” loop, getting the program to do something repeatedly	22
8	Date and times	24
9	Writing functions in R	25

```
# root dir
knitr::opts_knit$set(root.dir = "E:/cloudstor/IntroToR")
knitr::opts_chunk$set(echo = TRUE)
require(tidyverse)
require(lubridate)
```

1 Introduction

This is an introduction to R, written for ENVX3003. This work is based on earlier documents from the author and co-workers (in particular, Dasapta Erwin Irawan from ITBandung, Thomas Bishop and Floris van Ogtrop) and it also builds on many of the introduction to R literature on CRAN and elsewhere on the internet.

This course is not a complete introduction, and more in-depth knowledge on R and the use of R can be gained from many courses on-line and by basic practice.

This course covers three reading in files and plotting. It particularly uses the package **tidyverse** and Rstudio. More detail on how to use **tidyverse** is on this website.

We hope that this introduction offers sufficient depth to at least get you started with R and maybe later explore this in more depth yourself.

2 R as a modelling environment

The origins of R are in statistics, so this is what R does best. However, over time, it has proven to be a flexible language that can also be used quite effectively for programming and data science.

The key power of R is in using “scripts” and “notebooks” or “rmarkdown” files which are ways to record and document the code you are generating. In this course we will focus on using rmarkdown.

2.1 R and R Studio

2.1.1 Base R vs IDE

If R is the machine under the hood, then R Studio would be the dashboard, steering wheel, as well as the gas and brake paddles. People frequently refer to R as **base R** and R Studio is an Integrated Development Environment (IDE).

Are there other IDEs than R Studio? The answer is Yes. You could check out R Commander. Another interesting project is or Microsoft R Open, which offers a multithreaded version of R.

2.1.2 Running R online

Can we run R online? The answer is also Yes. R Studio offers a paid cloud service. You could try R fiddle for a limited range of code of package installation, CoCalc/Sage Math Cloud, Jupyter, and Code Ocean.

2.1.3 R is cross platform

R and R Studio are cross platform. So you could use R on all major operating systems (OS): Windows, Mac or Linux, so it's OK if you work with another person who doesn't use the same OS as you do. You just have to make sure that all parties have the same data and the same packages installed in the system, and the same code to run.

2.1.4 R components

In R, as in any other programming language, the two main components are the data and the codes. Using both, you could start an analysis and produce plots and tables as outputs. However in order to do some of the analyses, we will need **packages**. Packages are collections of functions that we can use in R.

The good thing about R is, there are *base functions*, that is commands that are included in the base R installation. This commands are progressing as you install newer versions of R. It's getting better and easier through time. But, because R is open source, users can develop their own scripts and functions or sets of functions. Sets of functions can be grouped as a *package*. So you would need to install the package first and load the package, before using the command or function inside that package. You would only need to install the package only once.

Generally we install a package from CRAN server using this basic code. Here **packageName** is the name of the package that you would like to install.

```
install.packages("packageName") # case sensitive and you need quotes  
library(packageName) # to load the package, case sensitive
```

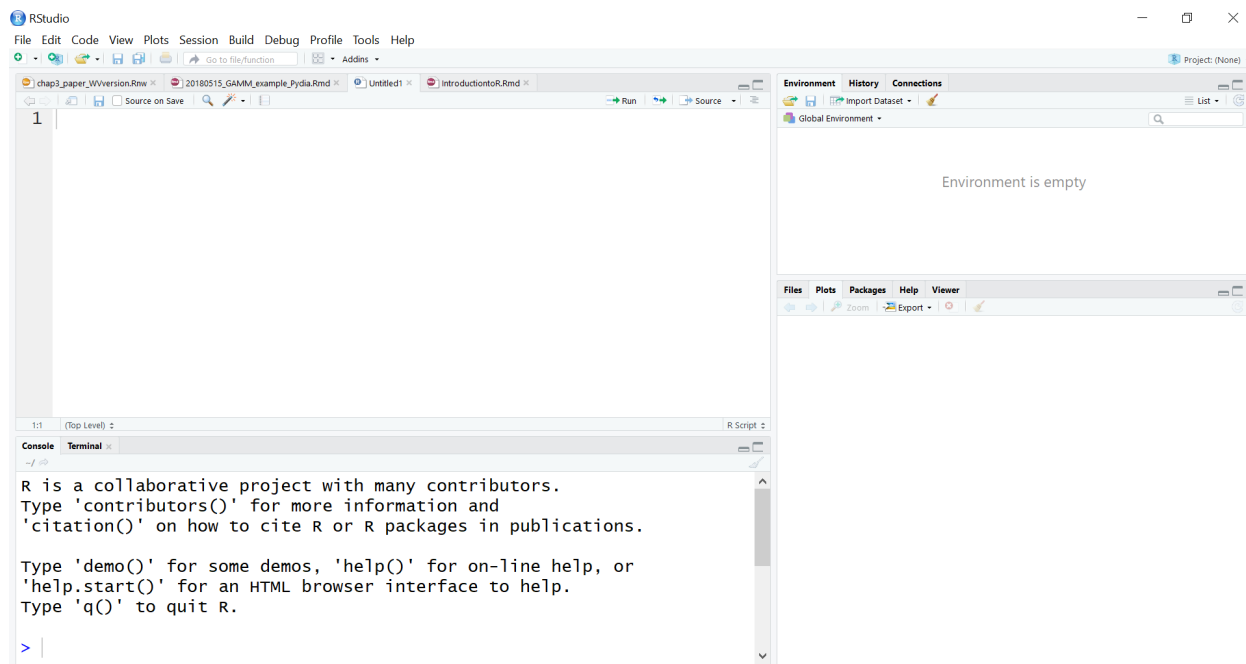


Figure 1: Four panels in R Studio

2.1.5 Navigation

If you use R Studio, you can see four panels (clock-wise): A *script* panel on top right, Environment, Files/folder/plots/packages, and console. You write your lines of code in the script panel then, click the *run* button (or select code and press CTRL+ENTER, or CMD+ENTER) to run and observe the progress of your code in the console panel. Find out in the console, if your code is running well or has a problem (error messages), or just a warning. Then you could see all the *objects* and loaded data components that relate to your code in the Environment panel.

3 Working folder structure

In R and in any other command line-based application, you would need to tell the program what your current folder location is and what the location of any data is. This is one of the most important steps in the process and many of the errors generated relate to not correctly define the directory. Usually we use the following folder structure:

- main project folder
- data: put your data here
- code or scripts: put your code here
- output: put your plots and tables here

However often we work with scripts, data, outputs in one folder, but use it as a process or intermediate folder. We usually sort out the components at the final stage of our work.

3.1 The working directory

Generally R works from a “working directory”. This is the directory on disk where it expects to find files or write files to. You can set this in Rstudio via the menu item “Session” → “Set working directory”, but you

can also set this using code in your script. Setting the working directory is useful when you want to access data in files on your computer or the network.

The basic function to use is `setwd("path/to/file")`. The thing to note is that in the path description you have to use “forward /” rather than the standard windows “backward”.

However, if you are using an Rmarkdown file, you are generally discouraged to use “setwd”, but here we can use another command which is part of the `knitr` package

3.1.0.1 Exercise

- Can you check your working folder/directory and what’s inside it?

in R:

`getwd()` # this tells you your current working directory
`dir()` # this gives you a list of the files in your currentt working directory
`?setwd` # this give you the help file on how to set your working directory in code

in Rmarkdown use `opts_knit$set(root.dir = "path/to/files")`

```
knitr::opts_knit$set(root.dir = "c:/users/rver4657/Desktop") ## windows users
knitr::opts_knit$set(root.dir = "~/Desktop") ## Mac users
dir() # this gives you a list of the files in your current working directory
```

4 BASIC R

4.1 R as a calculator

In its most basic form, R is a calculator

```
3*5
```

```
## [1] 15
```

```
50/100 + 0.1
```

```
## [1] 0.6
```

```
10 - 20
```

```
## [1] -10
```

4.2 Objects in R

The basic structure of R is based on objects, which are named. **R is case sensitive**, so keep this in mind. The main object we will use here is a *dataframe* or its modern variant the *tibble* in the package `tidyverse`.

All objects in R exist in the local R memory. So if you have a datafile, the first thing to do is to load it on your memory as an object that can be seen in the Environment panel. Thus, whatever you do with the object will not change your file, unless you save the object as a file.

R uses “<” to assign a value (or another object) to an object. You may find “=” means the same, but we don’t recommend it, because you also use “=” with different meaning in the command and parameter setting.

```
# assign
x <- 5
y <- 2
```

You can call up what is stored in the object (inspect) again by just typing its name:

```
x
```

```
## [1] 5
```

These objects will show up in the “Environment” window in Rstudio, or you can use `ls()` in the console to list the objects. The function `c()` can be used to stick things together into a vector. Redo the below commands in your own script.

```
# a vector
x <- c(1,2,5,7,8,15,3,12,11,19)
# another vector
y <- 1:10
# you have now two objects
ls()
```

```
## [1] "x" "y"
```

```
# you can add, multiply or subtract
z <- x + y
z
```

```
## [1]  2  4  8 11 13 21 10 20 20 29
```

```
zz <- x * y
zz
```

```
## [1]  1  4 15 28 40 90 21 96 99 190
```

```
zzz <- x - y
zzz
```

```
## [1]  0  0  2  3  3  9 -4  4  2  9
```

```
foo <- 0.5*x^2 - 3*x + 2
foo
```

```
## [1] -0.5 -2.0 -0.5  5.5 10.0 69.5 -2.5 38.0 29.5 125.5
```

- How many objects are now in your environment?

5 STATISTICAL ANALYSIS AND DATA MANIPULATION

Now it's time to look a bit further into more technical bits. How to manipulate data so we can perform some analyses on it to answer our research problem. There are, ofcourse, base R commands to do the job, but find it easier for us to use `tidyverse` package. This package is actually a combo of several packages written by the same author.

5.1 Packages to use

As highlighted, much of the power in R comes from the fact that it is open source and this means many people write new code and share this code. The formal way to do this is via “packages”, which, once checked and endorsed by the R community, appear in the CRAN repository as a **package**.

In this introduction we want to use some of the features in the package `tidyverse`. This package includes a series of other useful packages that you might need. The package is also the basis of the book *R for data science*

There are two components to using packages. The first is to make sure that the package is installed, for which we can use the functions `install.packages()`. Note that the name of the package is a *string* so needs to be between quotes `"`.

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

If the package is installed in your personal library, you will need to load the package in R using `require()` or `library()`. There are subtle differences between these two functions, but they are currently not that important. Check the help files.

```
require(tidyverse)
```

5.1.1 Exercise

- Can you load (and maybe first install) the package `lubridate`? This package is great for manipulating dates and times and works well with `tidyverse`.

5.2 A dataframe in tidyverse

A dataframe is a bit more complex. It is essentially a list, but presented as a table. Here is a simple demonstration of its power. We are using `data_frame()` from `tidyverse`.

```
Rainfall <- data_frame(City = c("Montevideo", "New York",  
                                "Amsterdam", "Sydney",  
                                "Moscow", "Hong Kong"),  
                       Rain_mm = c(950, 1174, 838, 1215,  
                                   707, 2400))
```

```
Rainfall
```

```
## # A tibble: 6 x 2  
##   City      Rain_mm  
##   <chr>      <dbl>  
## 1 Montevideo    950  
## 2 New York     1174  
## 3 Amsterdam     838  
## 4 Sydney      1215  
## 5 Moscow       707  
## 6 Hong Kong    2400
```

As you can see a dataframe (in this case a **tibble**) can mix character columns (`City`) and numeric columns (`Rain_mm`). Here I used `c()` to generate vectors which I put in the columns. In addition, the columns have names, which you can access using `colnames()` or `names()`:

```
colnames(Rainfall)
```

```
## [1] "City"      "Rain_mm"
```

```
names(Rainfall)
```

```
## [1] "City"      "Rain_mm"
```

Once you have a dataframe, you can access parts of the dataframe or manipulate the dataframe. Such as finding a column

```
# call a column
```

```
Rainfall$City
```

```
## [1] "Montevideo" "New York"    "Amsterdam"  "Sydney"     "Moscow"
## [6] "Hong Kong"
```

```
# or
```

```
Rainfall["City"]
```

```
## # A tibble: 6 x 1
##   City
##   <chr>
## 1 Montevideo
## 2 New York
## 3 Amsterdam
## 4 Sydney
## 5 Moscow
## 6 Hong Kong
```

```
# or
```

```
Rainfall[,1]
```

```
## # A tibble: 6 x 1
##   City
##   <chr>
## 1 Montevideo
## 2 New York
## 3 Amsterdam
## 4 Sydney
## 5 Moscow
## 6 Hong Kong
```

Subsetting rows is slightly different, you can still use numbers

```
# see the first two rows
```

```
Rainfall[1:2,]
```

```
## # A tibble: 2 x 2
##   City      Rain_mm
##   <chr>      <dbl>
## 1 Montevideo    950
## 2 New York     1174
```


5.2.1 filter()

But, because this is a tibble, we can use the power of tidyverse and the function `filter` to find rows

```
Rainfall %>%  
  filter(City=="Montevideo")
```

```
## # A tibble: 1 x 2  
##   City      Rain_mm  
##   <chr>      <dbl>  
## 1 Montevideo    950
```

OK, what is that `%>%` thing?? It is a further little symbol (apart from assign `<-`) in R that you need to know. It means “then”.

So to read the above code in words, it says:

- Take the Rainfall data, then, find the row where City equals “Montevideo”

Let’s expand this idea and show another `filter`:

```
# find a subset  
lots <- Rainfall %>%  
  filter(Rain_mm > 1000)  
lots
```

```
## # A tibble: 3 x 2  
##   City      Rain_mm  
##   <chr>      <dbl>  
## 1 New York    1174  
## 2 Sydney     1215  
## 3 Hong Kong   2400
```

Here I did two things, I did the filter using a comparison (and got multiple rows) and I “assigned” the result to a new object `lots`.

5.2.2 select()

Similar to `filter`, we can extract columns using `select`. So repeating the above example:

```
# call a column  
Rainfall %>%  
  select(City)
```

```
## # A tibble: 6 x 1  
##   City  
##   <chr>  
## 1 Montevideo  
## 2 New York  
## 3 Amsterdam  
## 4 Sydney  
## 5 Moscow  
## 6 Hong Kong
```

5.2.3 Exercise

Using the above examples, can you do the following?

- Extract the column with the rainfall values?
- Extract the row with the annual rainfall at Amsterdam?
- Which cities have rainfall below 1500 mm?

5.3 Reading data from different sources

There are a multitude of functions to read data from the disk into the R memory, I will demonstrate only one here, but more are given in the tidyverse book

Because a lot of data is stored in comma delimited txt files (such as Excel exports), using `read_csv()` is a good standard option.

Here I am reading in some monthly data from the Concordia station in the Uruguay river in Argentina. This data was originally downloaded from the Global River Discharge Database

```
UR_flow <- read_csv("Data/UruguayRiver_ConcordiaSt.csv")
```

```
## Parsed with column specification:
## cols(
##   Year = col_integer(),
##   Month = col_integer(),
##   Flow = col_integer()
## )
```

```
# check the first few lines (6 by default)
UR_flow
```

```
## # A tibble: 132 x 3
##   Year Month Flow
##   <int> <int> <int>
## 1  1969     1  7888
## 2  1969     2  5951
## 3  1969     3  4296
## 4  1969     4  4173
## 5  1969     5  4539
## 6  1969     6  4857
## 7  1969     7  4018
## 8  1969     8  3110
## 9  1969     9  2541
## 10 1969    10  2822
## # ... with 122 more rows
```

Previously you would have to save a specific program's data file, say in *xls* in to a pure text file such as *csv* or *txt*. However, there are now many packages that allow you to read a dataset directly from its binary format. There are many packages to do such task, readxl package is one of them. You could google your way of the most convenient package to use.

5.3.1 Exercise

- Can you read in the file: “Parana_CorrientesSt.csv” (supplied) from your data directory?

5.4 More data manipulation (using tidyverse)

5.4.1 Important commands

The following list is the important commands to remember:

- `select()` select columns
- `filter()` filter rows
- `arrange()` re-order or arrange rows (row sorting)
- `mutate()` create new columns
- `summarise()` summarise values
- `group_by()` allows for group operations in the “split-apply-combine” concept

Here is a short demonstration for `mutate` and `arrange` using the Rainfall data:

`mutate` and `select` to add columns and rearrange columns

```
# I can add a column of countries
Rainfall_new <- Rainfall %>%
  mutate(country = c("UY", "US", "NL", "AU", "RU", "CN")) %>%
  # and maybe a rainfall of the average monthly rainfall
  mutate(M_rain = Rain_mm/12)

# You can use select to reorder the columns and put country to the front
Rainfall_new <- Rainfall_new %>%
  select(country, everything())

# And if you would like to drop the M_rain column you can use
Rainfall_new %>%
  select(-M_rain)
```

```
## # A tibble: 6 x 3
##   country City      Rain_mm
##   <chr>   <chr>    <dbl>
## 1 UY      Montevideo  950
## 2 US      New York    1174
## 3 NL      Amsterdam   838
## 4 AU      Sydney      1215
## 5 RU      Moscow       707
## 6 CN      Hong Kong   2400
```

`arrange` to sort rows, `att desc()` do do this in decreasing order

```
Rainfall_new %>%
  arrange(desc(Rain_mm))
```

```
## # A tibble: 6 x 4
##   country City      Rain_mm M_rain
```

```
##   <chr>   <chr>       <dbl> <dbl>
## 1 CN     Hong Kong    2400 200
## 2 AU     Sydney       1215 101.
## 3 US     New York     1174 97.8
## 4 UY     Montevideo    950 79.2
## 5 NL     Amsterdam     838 69.8
## 6 RU     Moscow        707 58.9
```

Another useful function is `summarise()`, which allows you to apply a function over a data frame and particular across different factors. In tidyverse this is often combined with the function `group_by` to define how you would like to summarise.

Here is an example of summing the Uruguay river flow by year. Note that putting brackets around the statement makes it print out the result.

```
# aggregate to annual flow
(annual_flow <- UR_flow %>% #then
  group_by(Year=Year) %>% #then
  summarize(Sumflow = sum(Flow)))
```

```
## # A tibble: 11 x 2
##   Year Sumflow
##   <int> <int>
## 1 1969 53753
## 2 1970 52130
## 3 1971 66648
## 4 1972 99562
## 5 1973 103070
## 6 1974 47130
## 7 1975 68075
## 8 1976 53500
## 9 1977 73650
## 10 1978 41700
## 11 1979 61047
```

If you just wanted the overall mean and standard deviation of monthly flow, you could also ask for (and note how we drop `group_by`):

```
# Mean and sd\ monthly flow
UR_flow %>% #then
  summarize(Meanflow = mean(Flow),
            SdFlow = sd(Flow))
```

```
## # A tibble: 1 x 2
##   Meanflow SdFlow
##   <dbl> <dbl>
## 1 5457. 3492.
```

5.4.2 Exercise

- Can you calculate the standard deviation of the monthly flow by year using `summarize()`?

5.5 Different data set

Let's open this dataset. It's a water quality data in `csv` format.

Here we will introduce some groundwater chemical data from Semarang in Indonesia, kindly supplied by Dasapta Erwin Irawan from ITB

```
chemdata <- read_csv("data/semarang_chem.csv")
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   ID = col_character(),
##   Area = col_character(),
##   Year = col_integer(),
##   UTM_east = col_integer(),
##   UTM_north = col_integer(),
##   UTM_zone = col_character(),
##   Depth = col_integer(),
##   TDS = col_integer(),
##   EC = col_integer(),
##   Aq = col_character(),
##   Fac = col_character()
## )

## See spec(...) for full column specifications.
```

```
chemdata
```

```
## # A tibble: 58 x 24
##   ID      Area  Year  Lat  Long UTM_east UTM_north UTM_zone Depth  WL
##   <chr> <chr> <int> <dbl> <dbl>   <int>    <int> <chr>   <int> <dbl>
## 1 SB_1~ PT. ~ 1992 -6.96 110.  439936  9230800 49M      96 23.4
## 2 SB_2~ PT. ~ 1992 -6.98 110.  438500  9228100 49M      94 14.4
## 3 SB_2~ Obs.~ 1992 -6.96 110.  435500  9230150 49M     150 15.2
## 4 SB_2~ PT. ~ 1992 -6.98 110.  424150  9228400 49M      65 31.5
## 5 SB_2~ Dolo~ 1992 -6.97 110.  421950  9229400 49M      NA 19.8
## 6 SB_2~ Hote~ 1992 -6.99 110.  436950  9226950 49M      86  7.86
## 7 SB_3~ PT W~ 1992 -6.99 110.  427950  9226950 49M      76 44.6
## 8 SB_1~ PT. ~ 1992 -6.98 110.  436900  9228050 49M      NA 22.1
## 9 SB_2~ Tamb~ 1992 -6.98 110.  429800  9228600 49M      NA  4.64
## 10 SB_2~ Tamb~ 1992 -6.95 110.  423100  9231700 49M      80 11.6
## # ... with 48 more rows, and 14 more variables: Elev <dbl>, TDS <int>,
## #   ph <dbl>, EC <int>, K <dbl>, Ca <dbl>, Mg <dbl>, Na <dbl>, SO4 <dbl>,
## #   Cl <dbl>, HCO3 <dbl>, Bal <dbl>, Aq <chr>, Fac <chr>
```

A bit more on `select`. If you want multiple columns `Lat`, `Long` until `Depth`, you can simply use the `select()` function.

```
chemdata %>%
  select(Lat, Long:Depth)
```

```
## # A tibble: 58 x 6
##   Lat Long UTM_east UTM_north UTM_zone Depth
##   <dbl> <dbl>   <int>     <int> <chr>   <int>
## 1 -6.96 110.   439936   9230800 49M      96
## 2 -6.98 110.   438500   9228100 49M      94
## 3 -6.96 110.   435500   9230150 49M     150
## 4 -6.98 110.   424150   9228400 49M      65
## 5 -6.97 110.   421950   9229400 49M     NA
## 6 -6.99 110.   436950   9226950 49M      86
## 7 -6.99 110.   427950   9226950 49M      76
## 8 -6.98 110.   436900   9228050 49M     NA
## 9 -6.98 110.   429800   9228600 49M     NA
## 10 -6.95 110.   423100   9231700 49M      80
## # ... with 48 more rows
```

Or you want multiple columns Lat, Long until Depth, but you don't want UTM_zone. Again, you can use the `select()` function.

```
(chemdata %>%
  select(Lat, Long:Depth, -UTM_zone))
```

```
## # A tibble: 58 x 5
##   Lat Long UTM_east UTM_north Depth
##   <dbl> <dbl>   <int>     <int> <int>
## 1 -6.96 110.   439936   9230800    96
## 2 -6.98 110.   438500   9228100    94
## 3 -6.96 110.   435500   9230150   150
## 4 -6.98 110.   424150   9228400    65
## 5 -6.97 110.   421950   9229400   NA
## 6 -6.99 110.   436950   9226950    86
## 7 -6.99 110.   427950   9226950    76
## 8 -6.98 110.   436900   9228050   NA
## 9 -6.98 110.   429800   9228600   NA
## 10 -6.95 110.   423100   9231700    80
## # ... with 48 more rows
```

5.5.1 arrange()

Sorting out data by Aq and Fac. Use `arrange()` function.

```
chemdata %>%
  arrange(Aq, Fac)
```

```
## # A tibble: 58 x 24
##   ID Area Year Lat Long UTM_east UTM_north UTM_zone Depth WL
##   <chr> <chr> <int> <dbl> <dbl>   <int>     <int> <chr>   <int> <dbl>
## 1 SB_2~ PT. ~ 1992 -6.98 110.   424150   9228400 49M      65 31.5
## 2 SB_2~ Dolo~ 1992 -6.97 110.   421950   9229400 49M      NA 19.8
## 3 SB_3~ PT W~ 1992 -6.99 110.   427950   9226950 49M      76 44.6
## 4 SB_2~ PDAM~ 1992 -7.00 110.   431850   9226100 49M     100 56.6
## 5 SB_92 RS K~ 1992 -6.99 110.   434400   9227200 49M      75 6.8
## 6 SB_1~ Hote~ 1993 -7.00 110.   435763   9226415 49M     122 25.0
```

```
## 7 SB_2~ S. P~ 2003 -7.00 110. 432811 9225924 49M 152 5.97
## 8 SB_33 Es P~ 2003 -6.99 110. 429310 9227496 49M 90 22.5
## 9 SB_5~ Buki~ 2003 -6.99 110. 427050 9227750 49M 90 23.6
## 10 SP_3~ Obs.~ 2003 -6.99 110. 426474 9227278 49M NA 56.9
## # ... with 48 more rows, and 14 more variables: Elev <dbl>, TDS <int>,
## # ph <dbl>, EC <int>, K <dbl>, Ca <dbl>, Mg <dbl>, Na <dbl>, SO4 <dbl>,
## # Cl <dbl>, HCO3 <dbl>, Bal <dbl>, Aq <chr>, Fac <chr>
```

5.5.2 mutate()

Making new columns, for instance, calculating the ratio between Ca and Na. Use `mutate()` function

```
chemdata %>%
  mutate(ratio_Cana = Ca / Na)
```

```
## # A tibble: 58 x 25
##   ID   Area Year Lat Long UTM_east UTM_north UTM_zone Depth WL
##   <chr> <chr> <int> <dbl> <dbl> <int> <int> <chr> <int> <dbl>
## 1 SB_1~ PT. ~ 1992 -6.96 110. 439936 9230800 49M 96 23.4
## 2 SB_2~ PT. ~ 1992 -6.98 110. 438500 9228100 49M 94 14.4
## 3 SB_2~ Obs.~ 1992 -6.96 110. 435500 9230150 49M 150 15.2
## 4 SB_2~ PT. ~ 1992 -6.98 110. 424150 9228400 49M 65 31.5
## 5 SB_2~ Dolo~ 1992 -6.97 110. 421950 9229400 49M NA 19.8
## 6 SB_2~ Hote~ 1992 -6.99 110. 436950 9226950 49M 86 7.86
## 7 SB_3~ PT W~ 1992 -6.99 110. 427950 9226950 49M 76 44.6
## 8 SB_1~ PT. ~ 1992 -6.98 110. 436900 9228050 49M NA 22.1
## 9 SB_2~ Tamb~ 1992 -6.98 110. 429800 9228600 49M NA 4.64
## 10 SB_2~ Tamb~ 1992 -6.95 110. 423100 9231700 49M 80 11.6
## # ... with 48 more rows, and 15 more variables: Elev <dbl>, TDS <int>,
## # ph <dbl>, EC <int>, K <dbl>, Ca <dbl>, Mg <dbl>, Na <dbl>, SO4 <dbl>,
## # Cl <dbl>, HCO3 <dbl>, Bal <dbl>, Aq <chr>, Fac <chr>, ratio_Cana <dbl>
```

5.5.3 summarise()

Making a summary from your data. Use `summarise()` function.

```
chemdata %>%
  summarise(mean_TDS = mean(TDS),
            max_Cl = max(Cl),
            min_Cl = min(Cl),
            total = n())
```

```
## # A tibble: 1 x 4
##   mean_TDS max_Cl min_Cl total
##   <dbl> <dbl> <dbl> <int>
## 1 1041. 15753. 11.2 58
```

5.5.4 group_by()

Sorting out the data based on certain order. Use `group_by()` function.

```
chemdata%>%
  group_by(Aq) %>%
  summarise(mean_TDS = mean(TDS),
            max_Cl = max(Cl),
            min_Cl = min(Cl),
            total = n())
```

```
## # A tibble: 3 x 5
##   Aq                mean_TDS max_Cl min_Cl total
##   <chr>              <dbl>   <dbl> <dbl> <int>
## 1 Damar              371.     70    11.2    14
## 2 Garang             445.    146.    19.6    11
## 3 Quaternary marine 1523. 15753.    25     33
```

5.5.4.1 Exercise

- Can you calculate the `mean(Cl)` and `sd(Na)` for the dataset grouped by `Fac`?

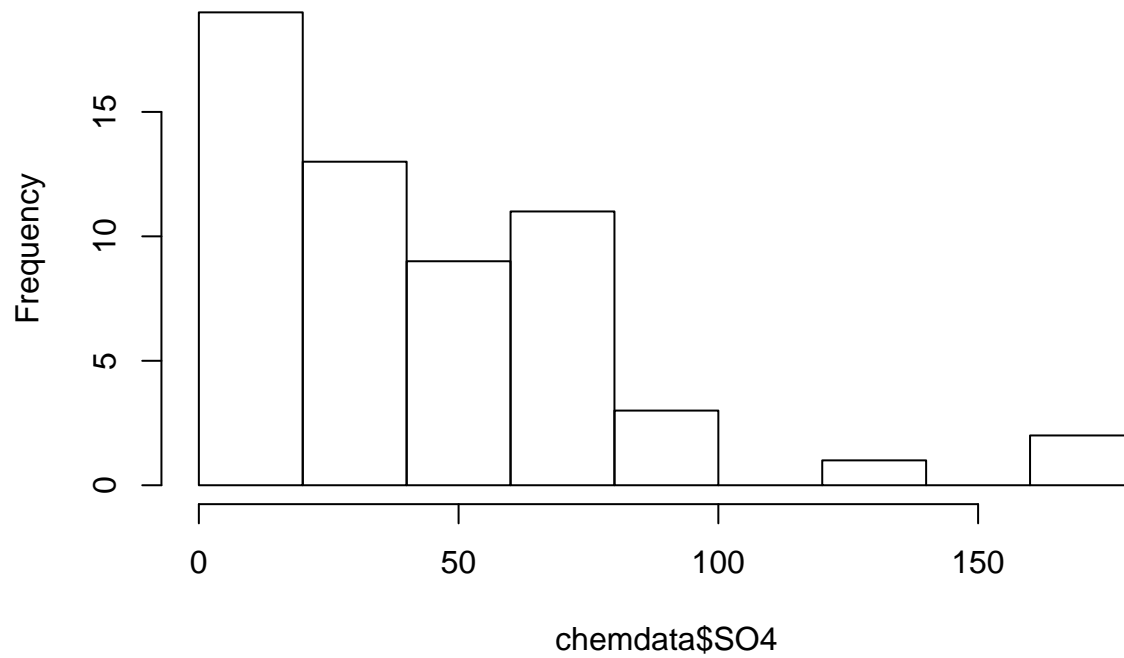
As we have indicated earlier, be sure to check out R for Data Science for more info about **tidyverse** and its use in data science.

6 PLOTTING

R is good at plotting. There are many ways to create a plot. So you just have to choose which one is the easiest for you. One way is using base R plotting engine. Like these plots.

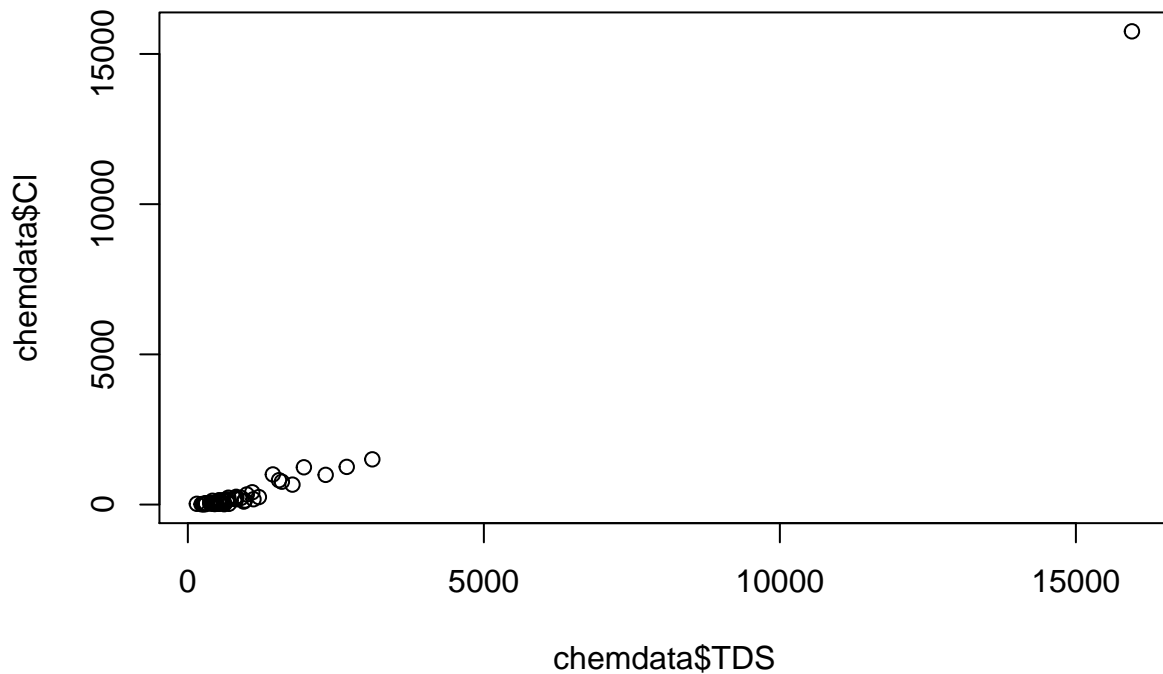
```
hist(chemdata$S04, main = "histogram of S04")
```


histogram of SO4



```
plot(chemdata$TDS, chemdata$Cl, main="scatter plot of TDS and Cl")
```

scatter plot of TDS and CI



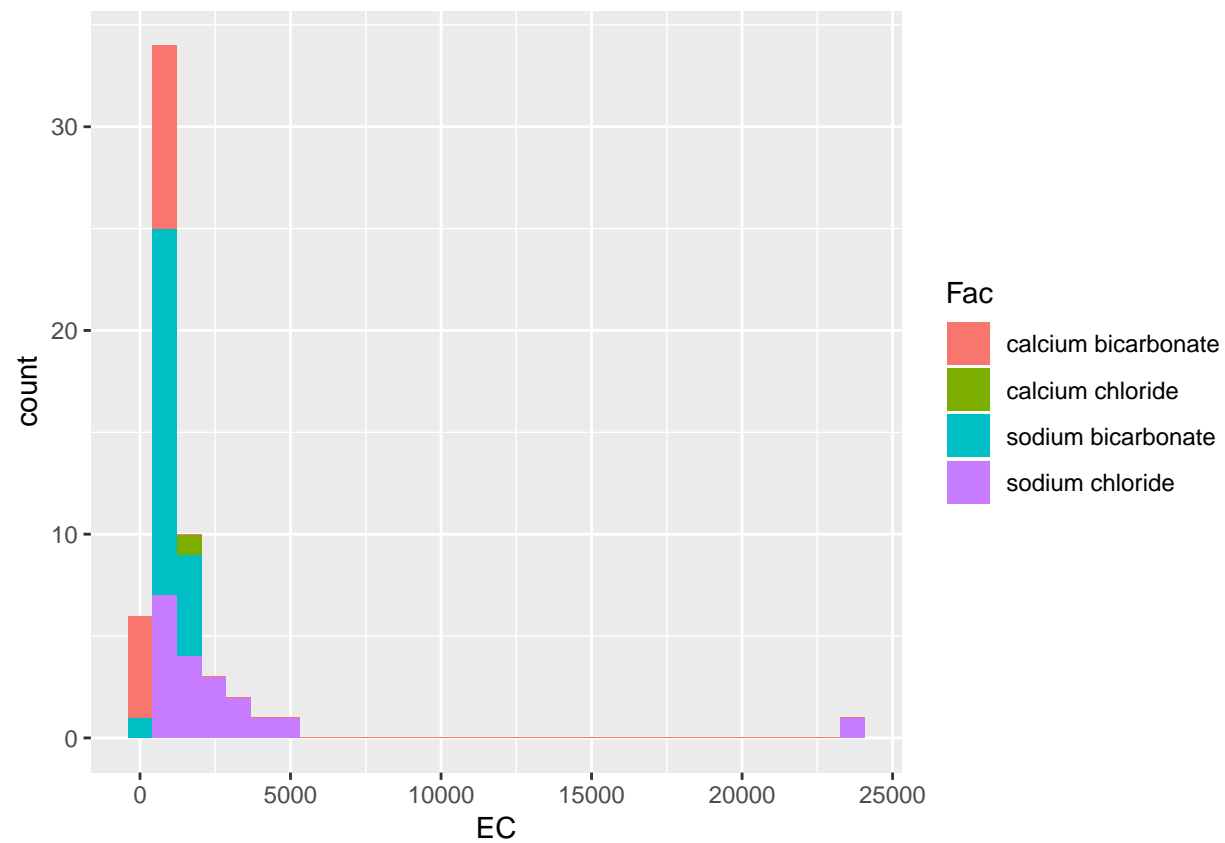
You could always tweak the plot to suits your needs. There are many resources about plotting in R, like:

- Producing Simple Graphs with R,
- Quick R.
- and more.

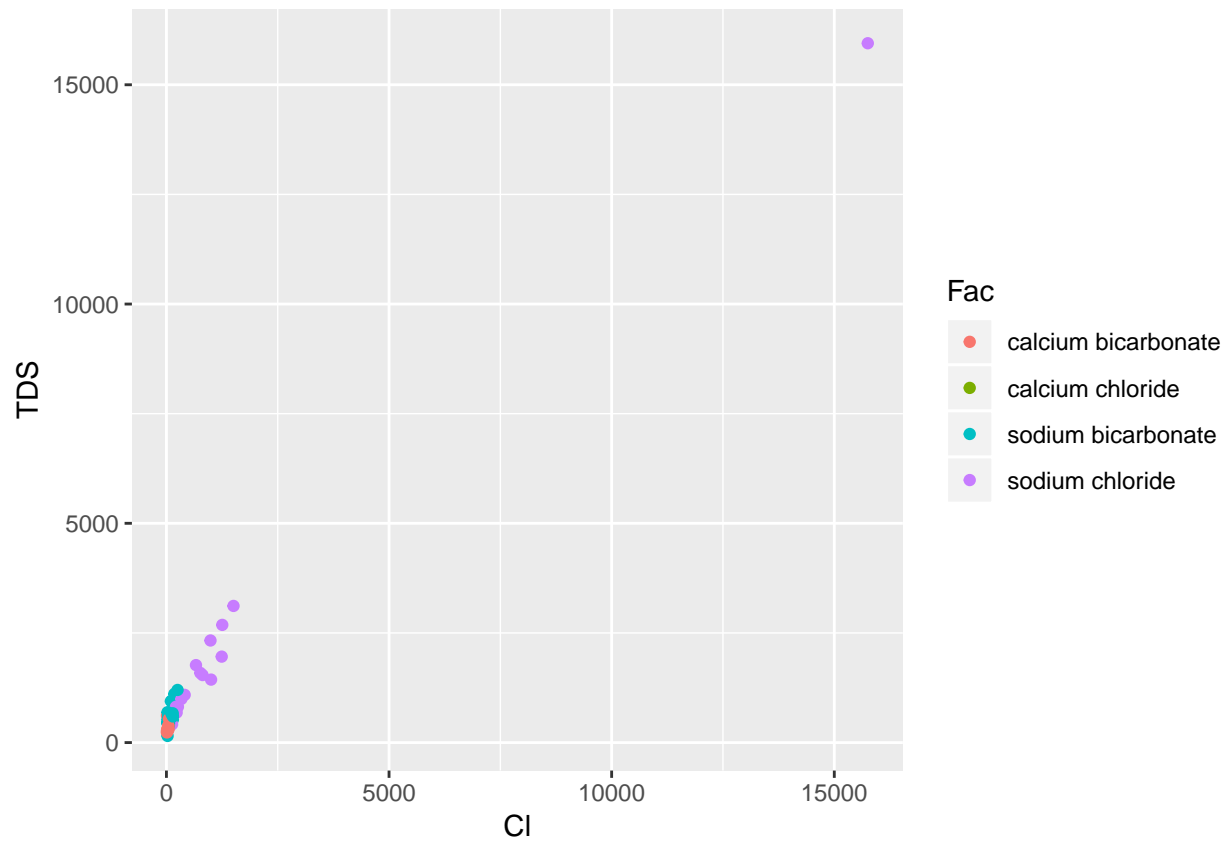
Nicer plots are made with `ggplot2` which is the plotting engine from `tidyverse`.

```
chemdata %>%  
  ggplot(aes(EC, fill = Fac)) + geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
chemdata %>%
  ggplot(aes(Cl, TDS, colour = Fac)) + geom_point()
```



In the second plot, note how colour is used to identify different facies.

6.1 Exercise

- try to make a plot between Ca vs Na using base R and ggplot2.
- try to make histogram for one parameter that you have in your dataset. Use ggplot2.
- can you tweak it by adding title to the plot and title to all axis. Hint: use `ggtitle()` and `labs()`.

7 Programming: if else and for loops

Loops are essential in programming. There are a range of different types of loops, but here I will only demonstrate “if else” and “for”. I have always been told that all other loops are just derivatives or short cuts. An “if else” loop allows you to program a switch in the code.

Basically it is used to evaluate an expression if the statement is TRUE and to evaluate another expression if the statement is FALSE:

if (comparison) do this, else do that

You might call the above line “pseudo code”. It is sometimes handy to first write something in pseudo code, basically you want to write in broad language what you want to happen.

You could also use only the “if” part, which means nothing happens if the statement is FALSE.

As an example I want to do: *if (a data frame has more than 10 rows) write data frame is LONG, else write data frame is SHORT.*

```
#We will first generate the data frame:
x <- UR_flow
# now wrote loop
if (nrow(x) > 10) {
  print("the dataframe is LONG")
} else {
  print("the dataframe is SHORT")
}
```

```
## [1] "the dataframe is LONG"
```

There are a few things to note here. I use the statement `nrow()` to check how many rows the data frame has. I use the statement `print()` to write something to the screen.

7.0.1 Exercise

- Try generating another data frame `x` and rerun the program, or change the program to get it to say “the dataframe is short”.

R also has an `ifelse()` command. This is a vectorized version of the `if` command - which means that it can be used on vectors of data - the command is applied to each element (or value) of the vector in turn. The `if` command only evaluates single values.

Using the `ifelse` command will return a vector of values, the same length as the longest argument in the expression.

Wherever possible, it is preferable to use the `ifelse` command rather than using the `if` command in combination with a loop - writing the program is more efficient and R evaluates vectorised functions more efficiently than it does loops. Here is an example which changes the program above.

```
x <- UR_flow
# add a column which identifies whether the flow < 5000
x[,4] <- ifelse(x[,3] > 5000, "large", "small")
# this creates a third column
tail(x,10)
```

```
## # A tibble: 10 x 4
##   Year Month Flow V4
```

```
##      <int> <int> <int> <chr>
## 1  1979      3  1699 small
## 2  1979      4  1650 small
## 3  1979      5  5052 large
## 4  1979      6  2619 small
## 5  1979      7  3776 small
## 6  1979      8  6134 large
## 7  1979      9  3275 small
## 8  1979     10 15692 large
## 9  1979     11 12244 large
## 10 1979     12  7380 large
```

I check in the second column of the data.frame whether the flows are greater than 5000 or not. I then write in the third column whether they are large or small numbers. A more complex (nested) ifelse version would be:

```
x[,5] <- ifelse(x[,3] > 2500, ifelse(x[,3] > 10000, "large", "intermediate"), "small")
tail(x, 10)
```

```
## # A tibble: 10 x 5
##   Year Month Flow V4    V5
##   <int> <int> <int> <chr> <chr>
## 1  1979      3  1699 small small
## 2  1979      4  1650 small small
## 3  1979      5  5052 large intermediate
## 4  1979      6  2619 small intermediate
## 5  1979      7  3776 small intermediate
## 6  1979      8  6134 large intermediate
## 7  1979      9  3275 small intermediate
## 8  1979     10 15692 large large
## 9  1979     11 12244 large large
## 10 1979     12  7380 large intermediate
```

You can try out some of your own versions of this

7.1 The “for” loop, getting the program to do something repeatedly

Loops are used to repeat a set of commands. Normally, there will be a variable which changes value in each successive loop through the commands. Reference to this changing value results in differences in output from successive iterations.

The for loop is used when the number of required iterations is known before the loop begins. It is used in the following way: *for (name in expression1) {expression2}*

- name is the name of the loop variable. Its value changes during each iteration, starting with the first value and ending with the last value in expression1.
- expression1 is a vector expression (often a sequence, such as 1:10).
- expression2 is a command or group of commands that are repeatedly evaluated. It usually contains references to name, which result in changes to the value of the expression as the value of name changes.

Here is the classic example of a loop

```
# Hello world
for (i in 1:5) {
print(paste(i, "hello world"))
}
```

```
## [1] "1 hello world"
## [1] "2 hello world"
## [1] "3 hello world"
## [1] "4 hello world"
## [1] "5 hello world"
```

Note the use of `paste()` to combine character vectors.
Here is another simple loop that tells you the first 5 values of the flow data.

```
for (i in 1:5) {
print(paste(UR_flow$Flow[i], "is the flow (ML/day)"))
}
```

```
## [1] "7888 is the flow (ML/day)"
## [1] "5951 is the flow (ML/day)"
## [1] "4296 is the flow (ML/day)"
## [1] "4173 is the flow (ML/day)"
## [1] "4539 is the flow (ML/day)"
```

```
# or more complex:
for (i in 1:5) {
print(paste("in Year", UR_flow$Year[i], "and month",
UR_flow$Month[i],
"the flow is", UR_flow$Flow[i], "(ML/day)"))
}
```

```
## [1] "in Year 1969 and month 1 the flow is 7888 (ML/day)"
## [1] "in Year 1969 and month 2 the flow is 5951 (ML/day)"
## [1] "in Year 1969 and month 3 the flow is 4296 (ML/day)"
## [1] "in Year 1969 and month 4 the flow is 4173 (ML/day)"
## [1] "in Year 1969 and month 5 the flow is 4539 (ML/day)"
```

You can also nest loops, that is, embed one loop into another. Here is an example that prints both the year and the flow using the column names in the dataframe.

```
for (i in 1:5) {
for (j in c(1,3)) {
print(paste(UR_flow[i,j], colnames(UR_flow)[j]))
}
}
```

```
## [1] "1969 Year"
## [1] "7888 Flow"
## [1] "1969 Year"
## [1] "5951 Flow"
## [1] "1969 Year"
```

```
## [1] "4296 Flow"
## [1] "1969 Year"
## [1] "4173 Flow"
## [1] "1969 Year"
## [1] "4539 Flow"
```

7.1.1 Exercise

- Write another program that includes a loop and a logical test

7.1.2 Comparison and Logical Operators

Comparison operators return a true or false value:

- `==` Equal to
- `>` Greater than
- `>=` Greater than or equal to
- `<` Less than
- `<=` Less than or equal to

Comparison operators can be combined with logical operators to describe more complex conditions.

Logical operators:

- `!` Not
- `|` or (used for vectors, with the `ifelse` command)
- `||` or (used for single values)
- `&` and (used for vectors, with the `ifelse` command)
- `&&` and (used for single values)

7.1.3 Exercise

Write a small program that uses comparison operators and a logical operator.

8 Date and times

This is a small demonstration of the package `lubridate`, which works well with `tidyverse` and allows you to convert times and dates with less effort. There is more description about this in this chapter of the `tidyverse` book.

```
#install.packages("lubridate") do this if you haven't done so.
require(lubridate)
UR_flow_d <- UR_flow %>%
  mutate(Dates = make_datetime(Year, Month))
```

Once we have this zoo data frame, it can be plotted quite easily with the basic plotting package.

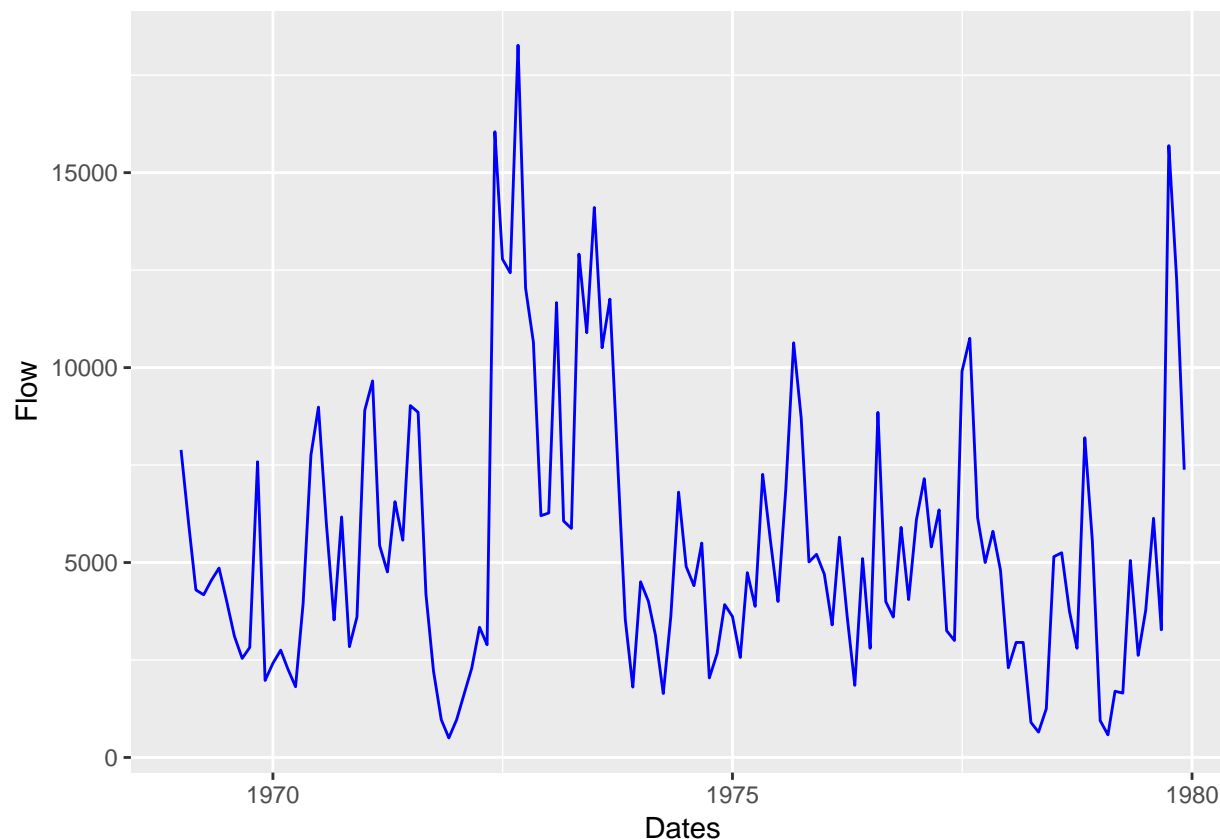


Figure 2: Demonstration of dates in plotting R

```
UR_flow_d %>%
  ggplot(aes(Dates,Flow)) + geom_line(colour="blue")
```

9 Writing functions in R

Until now you have used several functions in R that are part of packages or part of “core” R. However, another powerful element in R is the ability to write your own functions. There are two major advantages with writing functions:

1. They are easy to test, as they are contained. This is especially true if keep functions short.
2. They are short cuts and repeatable and therefore limit the possibility of typos.

Let’s go back to the “hello world” example that we used in a loop earlier. We can write the same example in a function.

The first thing to do is to decide which inputs we want the function to use to create the output. In this case I suggest we might want to change how many times the function produces output (which was 5 in the earlier example) and the actual output text, which was “hello world” in the original function.

The basic structure of a function is:

```
NameOfFunction <- function(input1, input2,...) {
  doSomething <- ....
```

```
return(doSomething)
}
```

Here is the hello world example:

```
HW <- function(n, outtext) {
  for (i in 1:n) {
    print(outtext)
  }
  # return("nothing")
}

# test
HW(5, "Hello World")
```

```
## [1] "Hello World"
## [1] "Hello World"
## [1] "Hello World"
## [1] "Hello World"
## [1] "Hello World"
```

```
# switch input by naming
HW(outtext = "I can switch the inputs", n = 3)
```

```
## [1] "I can switch the inputs"
## [1] "I can switch the inputs"
## [1] "I can switch the inputs"
```

Note that in this case the function produces output as part of its execution rather than returning an actual value (which is why I commented out the `return` statement). In the first example, you can see that you don't have to name the inputs if you keep the inputs in the same order as the defined function. R assumes that you mean `n = 5` and `outtext = "Hello world"`. In the second example I show that you can switch the inputs if you name them and that the function allows you to choose different inputs.

9.0.1 Exercise

- Can you write a function that calculates $y = a \cdot x + b$ for different values of `x`, `a` and `b`?
- Make the function return the output using `return()`

****END OF DOCUMENT****