

Review of the course “R for Data Science” Part 01(Talk 01~ 04)

By Haoran Nie @ HUST Life ST

Reference: R for Data Science

The book updated to 2nd ed. on July,2023, here's a [link](#) to the official website.

This work is licensed under CC BY-NC-SA 4.0 

Multi-omics data analysis and visualisation, #1

Talk 01

View the original slide through [this link](#).

View the original R markdown file of the slide through [this link](#).

This section has nothing to explain :)

R language basics, part 1

Talk 02

View the original slide through [this link](#).

View the original R markdown file of the slide through [this link](#).

Fundamental Data Type

The most basic data types include **numbers**, **logical symbols** and **strings** and are the basic building blocks of the other data types.

Simple Data Types

This includes vectors and matrices, both of which can contain multiple values of a certain basic data type, such as a **matrix** consisting of multiple numbers, a **vector** consisting of multiple strings, and so on. However, **they can only contain a single data type.**

```
1 c(100, 20, 30) ## Interger vector
2 c("String", "Array", "It's me".) ## String vector
3 c(TRUE, FALSE, TRUE, T, F) ## A logic vector
```

As shown above, arrays are usually defined with the function `c()`. In addition, a `vector` containing consecutive integers can be defined using the `:` operator.

Conversion between data types

1. Automatic Conversion

A `vector` can contain only one basic data type. Therefore, when defining arrays, if the input values are mixed, certain basic data types are automatically converted to other types to ensure consistency of the numeric types; this is called `coerce` in English, and has the meaning of forced

conversion. The priority of this conversion is:

- Logical types -> numeric types
- Logical Type -> String
- numeric type -> string

2. Manual switchover

In addition to the automatic conversion, you can manually convert the types of the elements in a vector:

- Checking the type of a variable `class()`
- Checking of classes `is.type()`
- Conversion of classes `as.type()`

Some special values in matrices

- `NA` (Not Available) missing values
- `NaN` (Not a Number) is meaningless
- `-Inf` Negative Infinity
- `Inf` Positive Infinity
- `NULL` Null

Some functions to determine these special values:

- `is.na()`
- `is.finite()`
- `is.infinite()`

Vectors and Arrays

Both are arrays. A `vector` is a one-dimensional array and a matrix is a two-dimensional array.

This means.

- There can be more dimensional arrays
- High-dimensional arrays, like `vector` and matrices, can contain only one basic data type.
- Higher dimensional arrays can be defined by the `array()` function.

Vector manipulation

```

1 dim(m);
2 nrow(m);
3 ncol(m);
4 range(m); ## Available when the content is numeric
5 summary(m); ## Can also be used in vector

```

Extra:

- Incorporation `ab = c(a, b)`
- Take part `ab[1]`
- Replacement of individual values `ab[1] = c`
- Replacing multiple values `ab[c(2, 3)] = c("Weihua", "Chen")`
- Naming elements and replace values `names(ab) = as.character(ab)`
- Reverse `rev(1:10)`
- Sort&order

```

1 lts = sample(LETTERS[1:20])
2 sort(lts)

```

- Fetch one line or multiple lines

```

1 # (There's already some data in workspace)
2
3 $ m
4 > (List the content of matrix "m")
5
6 $ m[1, ]
7 > (List the first row of matrix 'm')
8
9 $ m[1:2, ]
10 > (List the first two rows of matrix 'm')

```

You can also let the console to fetch multiple lines as the order you give.

```
1 m[c("row_B", "row_A")]
```

The console will output the contents of matrix "m" in the order of "row_B" and then "row_A".

- Fetch one column or multiple columns

As can be seen from the same principle, I only list codes here

```

1 m[, 1]
2 m[, c(1:2)]
3 m[, c("col_B", "col_A")]

```

- Fetch parts `m[1:2, 2:3]`
- Replacement

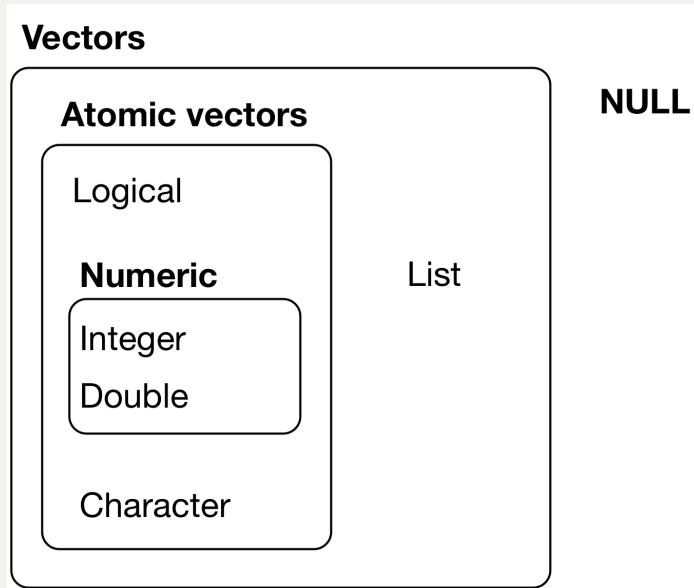
```

1 m[1, ] = c(10)
2 m[, "C"] = c(230, 140)
3 m[1:2,] = matrix( 1:6, nrow=2)
4 m[1, c("C", "B")] = matrix(110:111, nrow = 1)

```

- Transparent `t(m)`

The hierarchy of R's vector types



You can use function `typeof()` to know the type of a vector.

Here are some examples of other `is.xxx()` function:

```

1 is.null( NULL )
2 is.numeric( NA )
3 is.numeric( Inf );
4 is.list(); # This is a function which can take the place of
# "typeof()"
5 is.logical()
6 is.character()
7 is.vector();
8 # more ...
  
```

R language basics, part 2

Talk 03

View the original slide through [this link](#).

View the original R markdown file of the slide through [this link](#).

data.frame

What is a **data.frame**?

```
1 library(tidyverse);
2 library(kableExtra)
3 kbl(head(mpg),
4     booktabs = T)
```

Here's the result:

manufacturer	model	displ	year	cyl	trans	drv	cty	hwy	fl	class
audi	a4	1.8	1999	4	auto(l5)	f	18	29	p	compact
audi	a4	1.8	1999	4	manual(m5)	f	21	29	p	compact
audi	a4	2.0	2008	4	manual(m6)	f	20	31	p	compact
audi	a4	2.0	2008	4	auto(av)	f	21	30	p	compact
audi	a4	2.8	1999	6	auto(l5)	f	16	26	p	compact
audi	a4	2.8	1999	6	manual(m5)	f	18	26	p	compact

Usage of **head()** and **tail()**

- `head()` is a function to display the first rows of some data (vectors etc.)
- `tail()` is a function to display the last rows of some data (vectors etc.)

Components of `data.frame` and common functions

Components:

- Two-dimensional table
 - consists of different columns; each column is a vector, different columns can have different data types, but a column contains only one data type (`int`, `num`, `chr` ...)
 - Each column has the same length

Common functions:

```
1 nrow() # Show the number of rows  
2 ncol() # Show the number of columns  
3 dim() # Show the dimension
```

Structure of `data.frame` & `tibble`

```
str(mpg)
```

This command shows the structure of the tibble `mpg`:

Make a new `data.frame`

You can use the function `data.frame()` to make a new `data.frame`

```

1  data2 =
2    data.frame(
3      data = sample(1:100, 10),
4      group = sample(LETTERS[1:3], 10, replace = TRUE)
5      data2 = 0.1
6    )

```

How to add row(s)/col(s) to an existing `data.frame`

Create the "table header" first, then populate the `data.frame`

```

1  df2 =
2    data.frame(
3      x = character(),
4      y = integer(),
5      z = double() ,
6      stringsAsFactors = FALSE
7    )
8
9  df2 =
10   rbind(
11     df2,
12     data.frame(
13       x = "a",
14       y = 1L,
15       z = 2.2
16     )
17   )
18
19 df2 =
20   rbind(
21     df2,
22     data.frame(
23       x = "b",

```

```

24     y = 2,
25     z = 4.4
26   )
27 )

```

ATTENTION

- Use `rbind()` function to add rows, use `cbind()` function to add columns.
- Define the new line using `data.frame()` function, the "header" needs to be the same as the merged table.

You can also use these functions to bind several data.frames.

tibble

`tibble` is kind of similar to `data.frame`.

Make new tibble

`tibble` related functionality is provided by the `tibble` or `tidyverse` packages.

Almost all of the functions that you'll use in this book produce tibbles, as tibbles are one of the unifying features of the tidyverse. Most other R packages use regular data frames, so you might want to coerce a data frame to a tibble. You can do that with

`as_tibble()`:

```

1 as_tibble(iris)
2 #> # A tibble: 150 × 5
3 #>   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
4 #>     <dbl>      <dbl>      <dbl>      <dbl> <fct>
5 #> 1     5.1        3.5        1.4       0.2 setosa
6 #> 2     4.9        3          1.4       0.2 setosa
7 #> 3     4.7        3.2        1.3       0.2 setosa
8 #> 4     4.6        3.1        1.5       0.2 setosa
9 #> 5     5          3.6        1.4       0.2 setosa
10 #> 6    5.4        3.9        1.7       0.4 setosa
11 #> # i 144 more rows

```

Another way to create a tibble is with `tribble()`, short for **t**ransposed **t**ibble.

`tribble()` is customised for data entry in code: column headings are defined by formulas (i.e. they start with `~`), and entries are separated by commas. This makes it possible to lay out small amounts of data in easy to read form.

```

1 tribble(
2   ~x, ~y, ~z,
3   #--|--|----
4   "a", 2, 3.6,
5   "b", 1, 8.5
6 )
7 #> # A tibble: 2 × 3
8 #>   x         y     z
9 #>   <chr> <dbl> <dbl>
10 #> 1 a         2     3.6
11 #> 2 b         1     8.5

```

- `add_row()`
- `add_column()`

Manipulate the tibble

See “Manipulate the `data.frame`”

tibble to data.frame

- `as.data.frame()`
- `as_tibble()`

e.g.

```
1 library(tibble)
2 as.data.frame(head(as_tibble(iris)))
```

Differences between tibble and data.frame

Tibble evaluates columns sequentially

```
1 rm(x,y) # Delete possible x, y
2 tibble(x = 1:5, y = x^2); # You can do this with tibble
3 data.frame(x = 1:5, y = x ^ 2); # But data.frame doesn't work.
```

`data.frame` causes trouble when fetching `subset` operations

```
1 df1 =
2   data.frame(x = 1:3, y = 3:1)
3 class(df1[, 1:2])
4
5 #> [1] "data.frame"
6
7 # Subset operation :takes a column and expects a data.frame ()
8 class(df1[, 1]) # The result is a vector ...
9
10 #> [1] "integer"
11
```

```

12 ## Tibble doesn't.
13 df2 =
14   tibble(x = 1:3, y = 3:1)
15 class(df2[, 1]) ## Tibble forever
16
17 #> [1] "tbl_df" "tbl" "data.frame"

```

`tibble` allows controlled data type conversion

There's no proper example here.

:_(

Recycling

```

1 data.frame(a = 1:6, b = LETTERS[1:2]) # data.frame CAN!!!

```

OUTPUT

```

1 #   a b
2 # 1 1 A
3 # 2 2 B
4 # 3 3 A
5 # 4 4 B
6 # 5 5 A
7 # 6 6 B

```

```

1 tibble(a = 1:6, b = LETTERS[1:2]); ## But tibble CAN'T!!!

```

OUTPUT

```

1 # Error:
2 # ! Tibble columns must have compatible sizes. ## * Size 6:
3 # * Size 2: Column `b`.
4 # i Only values of size one are recycled.

```

ATTENTION!

The recycling of `tibble` is limited to lengths of 1 or equal; `data.frame` is just divisible.

`data.frame` will do partial matching, while `tibble` will NEVER do it.

```

1 df = data.frame(abc = 1)
2 df$ab; # Unwanted result ...
3
4 df2 = tibble(abc = 1)
5 df2$a; # Produce a warning and return NULL

```

OUTPUT

```

1 # Warning: Unknown or uninitialized column: `a`.
2 # NULL

```

Advanced tips for using `data.frame` and `tibble`

- `attach()`
- `detach()`
- `with()`
- `within()`

Following is the introduction (Produced by ChatGPT)

These functions—`attach()`, `detach()`, `with()`, and `within()`—are incredibly useful when working with data frames or tibbles in R, aiding in smoother workflows and code readability. Here's a breakdown of their functionality:

`attach()` and `detach()`

- **Purpose:** These functions allow you to temporarily attach a data frame to the search path, making its columns directly accessible by their names.
- **Usage:**
 - `attach(df)` attaches the specified data frame `df`.
 - `detach(df)` detaches the specified data frame `df`.
- **Example:**

```

1 data(mtcars) # Loading a sample dataset
2 attach(mtcars) # Attaching mtcars
3
4 # Now, columns can be accessed directly
5 summary(mpg)
6 mean(mpg)
7
8 detach(mtcars) # Detaching mtcars

```

- **Note:** While convenient, using `attach()` can sometimes lead to confusion or unintended consequences, such as masking variables in your environment. It's often recommended to avoid using `attach()` due to potential side effects.

`with()`

- **Purpose:** `with()` allows you to execute expressions in an environment where the data frame's columns can be referenced without using `$`.
- **Usage:**
 - `with(data, expr)` evaluates `expr` in the context of the specified data frame `data`.

- **Example:**

```

1 data(mtcars) # Loading a sample dataset
2
3 with(mtcars, {
4   mean(mpg)
5   summary(cyl)
6 })

```

- **Advantage:** It helps avoid repetitive use of the data frame name while working with its columns.

`within()`

- **Purpose:** Similar to `with()`, `within()` allows modification of a data frame by evaluating expressions within it.
- **Usage:**
 - `within(data, expr)` modifies `data` according to `expr` and returns the modified data frame.
- **Example:**

```

1 data(mtcars) # Loading a sample dataset
2
3 modified_mtcars <- within(mtcars, {
4   mpg_square <- mpg^2
5   hp_doubled <- hp * 2
6 })
7 head(modified_mtcars) # Checking the modified data
frame

```

- **Advantage:** `within()` is useful when you want to create or modify columns within the data frame without having to repeatedly refer to the data frame name.

Remember, while these functions can streamline your code, it's crucial to use them judiciously to avoid unexpected behavior or cluttering your global environment.

Here's the console print output for the examples provided earlier:

```

1 # Using attach() and detach()
2 data(mtcars) # Loading a sample dataset
3 attach(mtcars) # Attaching mtcars
4
5 # Now, columns can be accessed directly
6 summary(mpg)
7 # Output:
8 #   Min. 1st Qu. Median   Mean 3rd Qu.   Max.
9 # 10.40 15.43 19.20 20.09 22.80 33.90
10
11 mean(mpg)
12 # Output:
13 # [1] 20.09062
14
15 detach(mtcars) # Detaching mtcars
16
17 # Using with()
18 data(mtcars) # Loading a sample dataset
19
20 with(mtcars, {
21   mean(mpg)
22   # Output:
23   # [1] 20.09062
24
25   summary(cyl)
26   # Output:
27   #   Min. 1st Qu. Median   Mean 3rd Qu.   Max.
28   # 4.00 4.00 6.00 6.188 8.00 8.00
29 })
30
31 # Using within()
32 data(mtcars) # Loading a sample dataset
33
34 modified_mtcars <- within(mtcars, {
35   mpg_square <- mpg^2
36   hp_doubled <- hp * 2
37 })
```

```

38 head(modified_mtcars) # Checking the modified data frame
39 # Output:
40 #          mpg cyl disp  hp drat    wt  qsec vs am
41 gear carb mpg_square hp_doubled
42 # Mazda RX4           21.0   6 160 110 3.90 2.620 16.46 0  1
43   4      441.00        220
44 # Mazda RX4 Wag       21.0   6 160 110 3.90 2.875 17.02 0  1
45   4      441.00        220
46 # Datsun 710          22.8   4 108  93 3.85 2.320 18.61 1  1
47   4      519.84        186
48 # Hornet 4 Drive      21.4   6 258 110 3.08 3.215 19.44 1  0
49   3      457.96        220
50 # Hornet Sportabout   18.7   8 360 175 3.15 3.440 17.02 0  0
51   3      349.69        350
52 # Valiant              18.1   6 225 105 2.76 3.460 20.22 1  0
53   3      327.61        210

```

File IO

Read from files

Using functions from the `readr` package

```

1 # readr is part of tidyverse
2 library(tidyverse) # or alternatively
3 library(readr)

```

Some available functions:

- `read_csv()`: comma separated (CSV) files
- `read_tsv()`: tab separated files
- `read_delim()`: general delimited files
- `read_fwf()`: fixed width files
- `read_table()`: tabular files where columns are separated by white-space.
- `read_log()`: web log files

Full documentation of the package is available through this [link](#).

Usage

- Read with predefined column types

```

1 myiris2 =
2   read_csv("../data/talk03/iris.csv",
3           col_types =
4             cols(
5               Sepal.Length = col_double(),
6               Sepal.Width = col_double(),
7               Petal.Length = col_double(),
8               Petal.Width = col_double(),
9               Species = col_character()
10              )
11            )

```

- To read from other formats, you can try the following packages:

Similar to Python

- `haven` - SPSS, Stata, and SAS files
- `readxl` - excel files (.xls and .xlsx) DBI - databases
- `jsonlite` - json
- `xml2` - XML
- `httr` - Web APIs
- `rvest` - HTML (Web Scraping)

Write to files

Use the following functions to write object(s) to external files:

Default parameters are listed.

More related documents can be found in this [link](#).

- Comma delimited file:

```

1  write_csv(
2    x,
3    path,
4    na = "NA",
5    append = FALSE,
6    col_names = !append
7  )

```

- File with arbitrary delimiter:

```

1  write_delim(
2    x,
3    path,
4    delim = " ",
5    na = "NA",
6    append = FALSE,
7    col_names = !append
8  )

```

- CSV for excel:

```

1  write_excel_csv(
2    x,
3    path,
4    na = "NA",
5    append = FALSE,
6    col_names = !append
7  )

```

- String to file:

```

1  write_file(
2    x,
3    path,
4    append = FALSE
5  )

```

- String vector to file, one element per line:

```

1  write_lines(
2    x,
3    path,
4    na = "NA",
5    append = FALSE
6  )

```

- Object to RDS file:

```

1  write_rds(
2    x,
3    path,
4    compress =
5      c(
6        "none",
7        "gz",
8        "bz2",
9        "xz"
10       ),
11     ...
12   )

```

- Tab delimited files:

```
1 write_tsv(  
2     x,  
3     path,  
4     na = "NA",  
5     append = FALSE,  
6     col_names = !append  
7 )
```

R language basics, part 3: factor

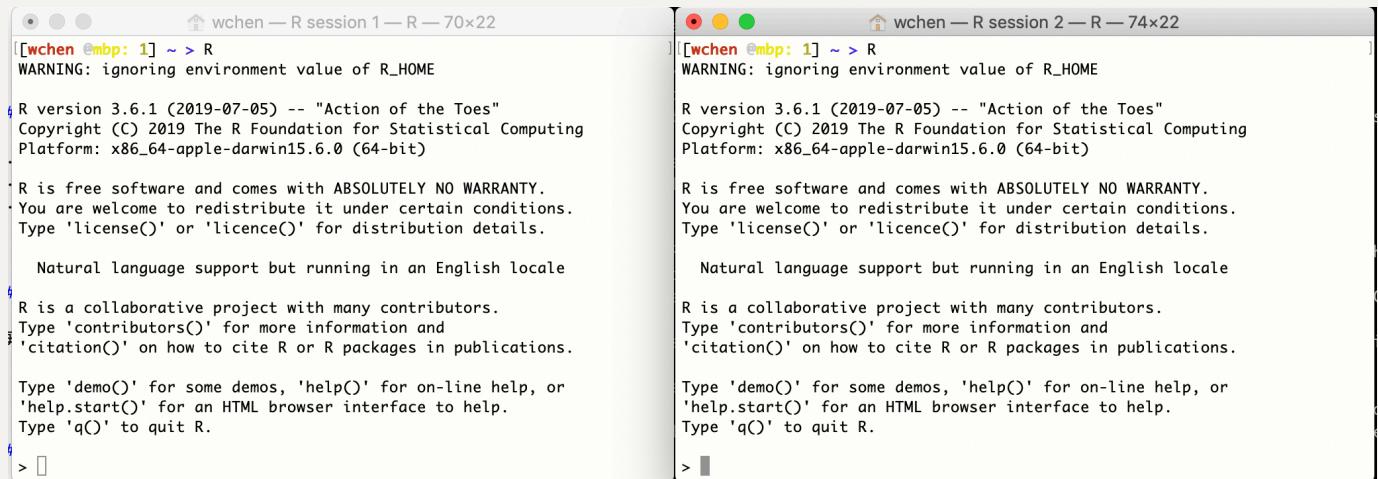
Talk 04

View the original slide through [this link](#).

View the original R markdown file of the slide through [this link](#).

IO and working environment management

Each R session is a separate **work space** containing its own data, variables, and operation history.



```
wchen — R session 1 — R — 70x22
[wchen @mbp: 1] ~ > R
WARNING: ignoring environment value of R_HOME

R version 3.6.1 (2019-07-05) -- "Action of the Toes"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (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.

> 
```



```
wchen — R session 2 — R — 74x22
[wchen @mbp: 1] ~ > R
WARNING: ignoring environment value of R_HOME

R version 3.6.1 (2019-07-05) -- "Action of the Toes"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (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.

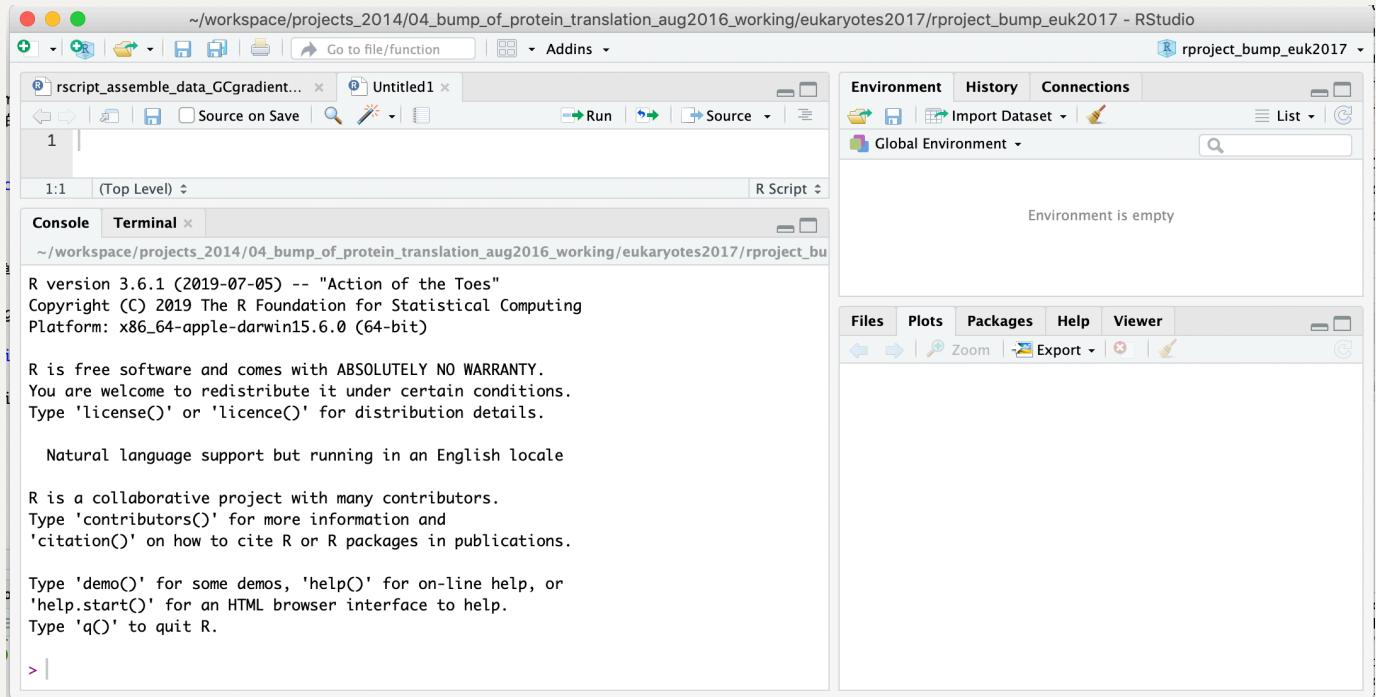
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> 
```

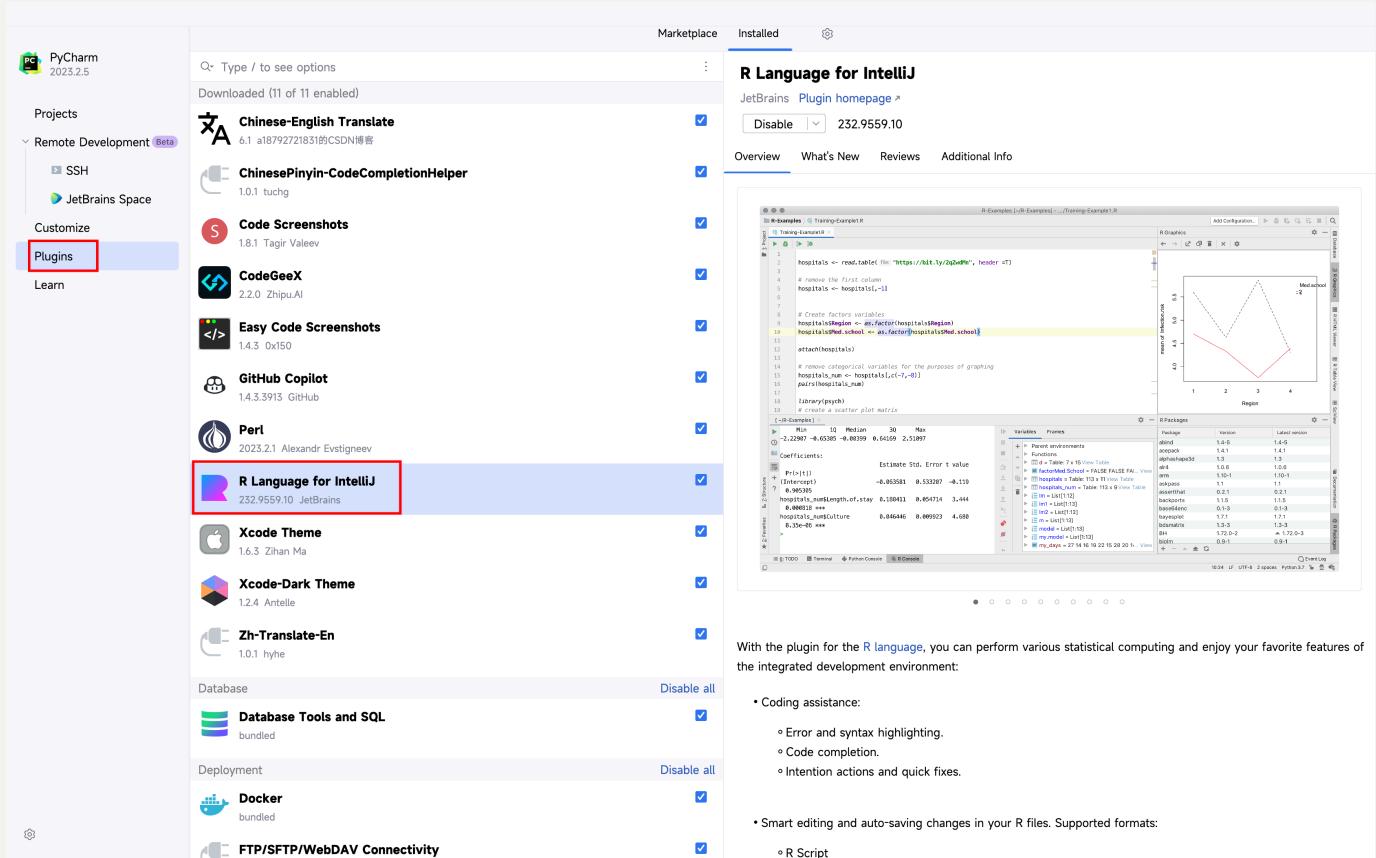
Each RStudio session is automatically associated with a R session

Not only RStudio, PyCharm or VSCode also support R session.

However, I'm keen on coding with PyCharm but not RStudio, for its wonderful Plug-in Environment, which can let me use plug-ins such as Code GeeX by Zhipu AI (a company founded by some student in KEG team in Tsinghua University) or GitHub Copilot by GitHub to let the coding process more quickly, for the instruction from GPTs.



If you want to coding with R using PyCharm or other JetBrains IDE (i.e. IntelliJ, CLion, etc.), remember to install the *R Language Plug-in*



For instruction how to get FREE Student Liscence of GitHub Pro, GitHub Copilot and JetBrains Products and their benefits, see their official website:

- GitHub Global Campus

Make sure you don't use VPNs and use your phone to log in and apply (HUST Campus Network is recommended), give "Precise Location" permission to your browser. You may use your "[Student Number]@hust.edu.cn" mail to verify your identity as a student studying in HUST.

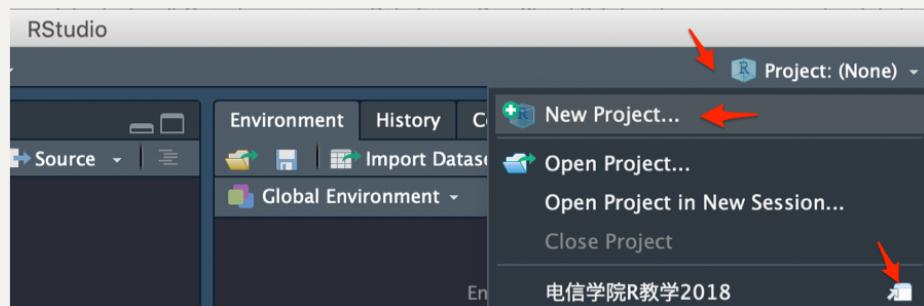
- JetBrains Products

Because our email addresses ending with "@hust.edu.cn" are banned due to misuse, you should apply for an online verification report on [CHSI](#) (press the link to visit the website), instructions [here](#).

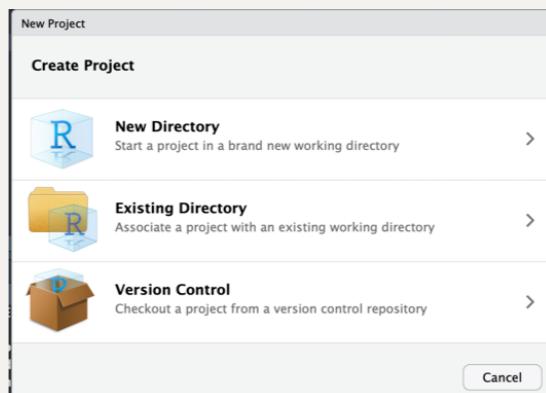
Start a new RStudio session by creating a new project

To start a new session in PyCharm, simply press the bottom corner and select a new session.

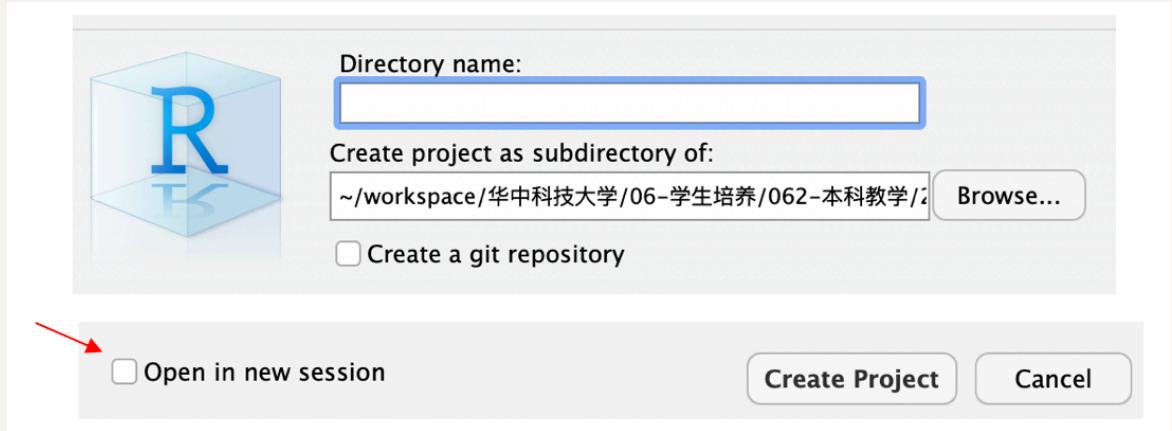
- Click the Project button in the upper right corner and select New Project in the pop-up menu ...



- Select: New directory -> New Project in the popup window



- Enter a new directory name, choose its mother directory ...



Working Space

Current workspace, including all loaded data, packets and homebrew functions.

Variables can be managed with the following code:

```
1 ls() # Show all the variables in current workspace/session
2 rm(x) # Remove a variable
3 rm(list = ls()) # Remove ALL variables in current
      workspace/session
```

Variables in working space in RStudio

The "Environment" window in the upper right corner of RStudio shows all the variables of the current workspace.

Data	
▶ aq	153 obs. of 7 variables
▶ dat	10 obs. of 3 variables
▶ dat2	10 obs. of 3 variables
▶ df	1 obs. of 1 variable
▶ df2	1 obs. of 1 variable
▶ l	List of 26
m	num [1:2, 1:3] 1 2 111 103 110 101
▶ myiris	150 obs. of 5 variables
▶ myiris2	150 obs. of 5 variables

Values	
a	int [1:3] 1 2 3
ab	Named chr [1:6] "1" "—" "≡" "ah" "bo" "C"
b	chr [1:3] "A" "B" "C"
lts	chr [1:20] "P" "N" "I" "S" "Q" "D" "C" "A" "J" "M" "E" "F" "H"...

Save and restore work space

```

1 # Save all loaded variables into an external .RData file
2 save.image(file = "prj_r_for_bioinformatics_aug3_2019.RData")
3 # Restore (load) saved work space
4 load(file = "prj_r_for_bioinformatics_aug3_2019.RData")

```

Notes:

- Existing variables will be kept, however, those with the same names will be replaced by loaded variables
- Please consider using `rm(list=ls())` to remove all existing variables to have a clean start
- You may need to reload all the packages

Save selected variables

Sometimes you need to transfer processed data to a collaborator ...

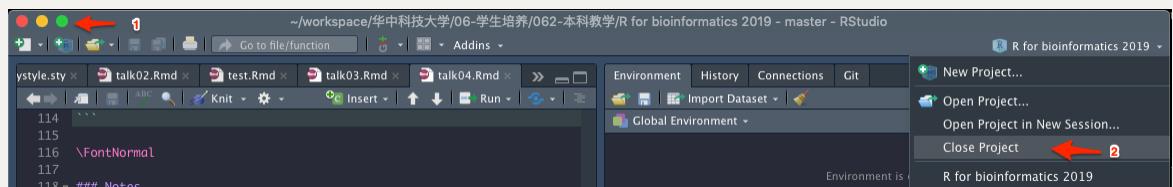
```

1 # Save selected variables to external
2 save(
3   city,
4   country,
5   file="1.RData"
6 )
7 # You can specify directory name
8 load("1.RData")

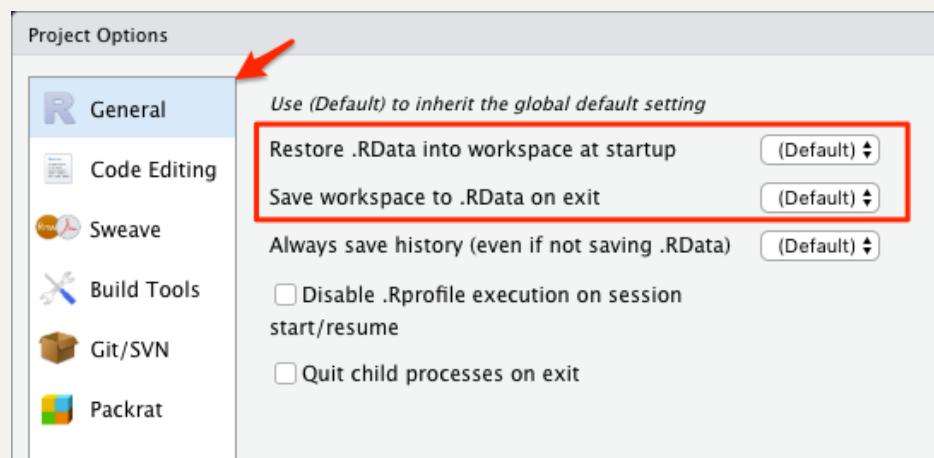
```

Close and (re)open a project

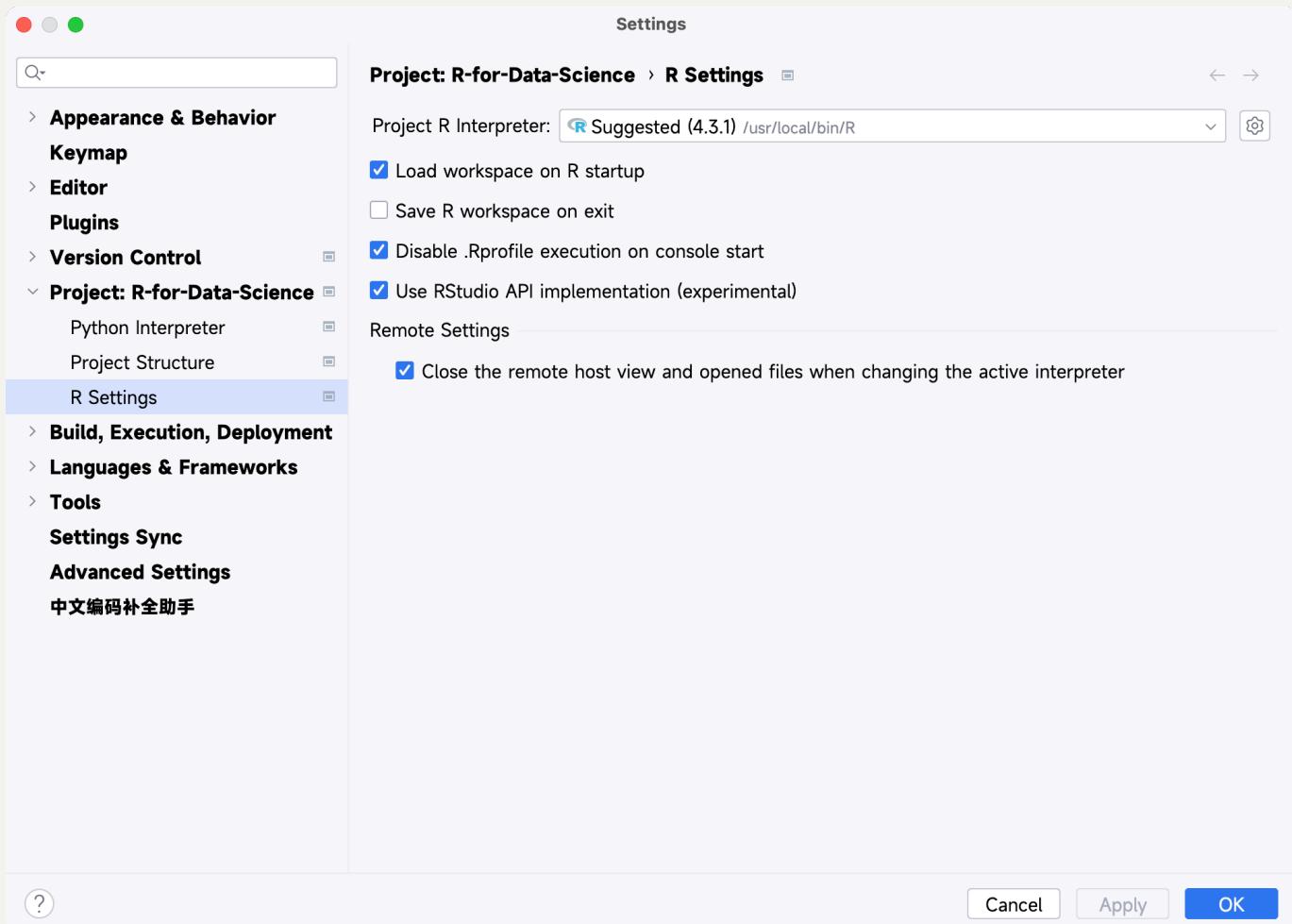
- To close a project



- In RStudio and similar IDEs, there are some preferences to choose



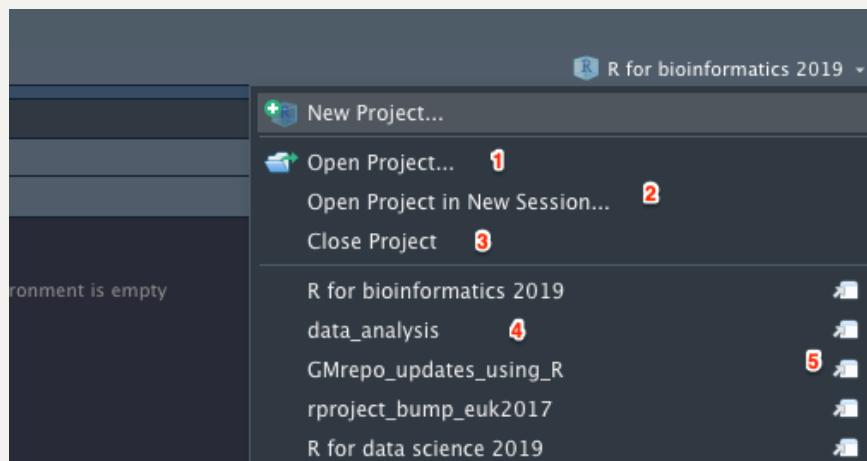
The UI in PyCharm



Notes:

- Save on exit
- Load on opening
- When the data is large, the loading time may be too long ...

Open a project



When in PyCharm, simply drag the working directory to its main window, remember to trust the project.

Factors

Factor is a data structure used for fields that takes only predefined, finite number of values (categorical data).

It will limit the selection of input data.

Play around with `levels()`

Here are instructions of modifying factor levels

Based on the textbook

The levels are terse and inconsistent. Let's tweak them to be longer and use a parallel construction. Like most rename and recoding functions in the tidyverse, the new values go on the left and the old values go on the right:

```

1  load(gss_cat)
2
3  mutate(
4    partyid = fct_recode(partyid,
5      "Republican, strong"     = "Strong republican",
6      "Republican, weak"      = "Not str republican",
7      "Independent, near rep" = "Ind,near rep",
8      "Independent, near dem" = "Ind,near dem",
9      "Democrat, weak"        = "Not str democrat",
10     "Democrat, strong"      = "Strong democrat"
11   )
12 )
13
14 count(partyid)
15
16 #> # A tibble: 10 × 2
17 #>   partyid          n

```

```

18 #> <fct>                <int>
19 #> 1 No answer            154
20 #> 2 Don't know          1
21 #> 3 Other party         393
22 #> 4 Republican, strong 2314
23 #> 5 Republican, weak   3032
24 #> 6 Independent, near rep 1791
25 #> # i 4 more rows

```

Use this technique with care: if you group together categories that are truly different you will end up with misleading results.

The order of the `levels` determines the sorting order.

Use factor to clean data

Usage of `fct_xxx()` functions.

Suppose I have a set of gender data that is written in a very irregular way:

```

1 gender =
2   c("f", "m ", "male ", "male", "female", "FEMALE", "Male", "f",
3     "m" )
4
5 gender_fct =
6   as.factor(gender)
7 fct_count(gender_fct)

```

The output looks like this:

```
> fct_count(gender_fct)
# A tibble: 8 × 2
   f          n
   <fct>    <int>
1 "f"        2
2 "female"   1
3 "FEMALE"   1
4 "m"        1
5 "m "       1
6 "male"     1
7 "Male"     1
8 "male "    1
```

Now I request to replace with Female, Male.

```
1 gender_fct =
2   fct_collapse(
3     gender,
4     Female = c("f", "female", "FEMALE"),
5     Male = c("m ", "m", "male ", "male", "Male")
6   )
7
8 fct_count(gender_fct)
```

```
# A tibble: 2 × 2
   f          n
   <fct>    <int>
1 Female     4
2 Male       5
```

You can also use `fct_relabel()` to do the same thing

```

1 fct_relabel(
2   gender,
3   ~ ifelse(
4     tolower(
5       substring(., 1, 1)) == "f",
6       "Female",
7       "Male"
8     )
9 )

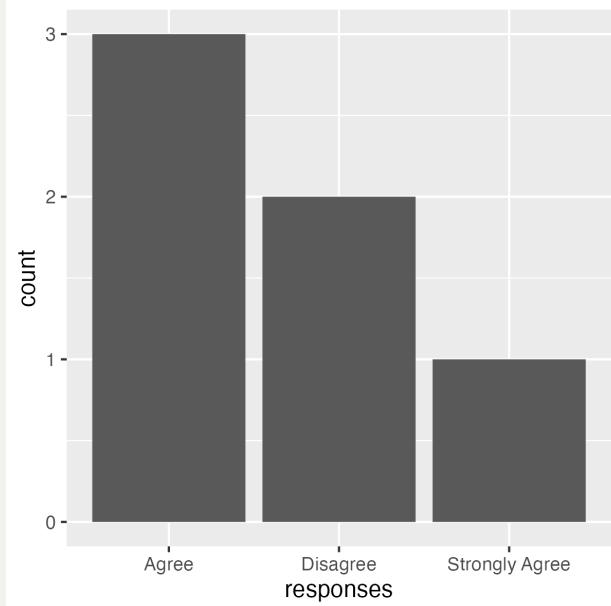
```

Usage of factors in drawing plots

```

1 library(ggplot2)
2
3 responses =
4   factor(
5     c("Agree", "Agree", "Strongly Agree", "Disagree",
6       "Disagree", "Agree")
7     )
8
9 response_barplot =
10   ggplot(
11     data = data.frame(responses),
12     aes(x = responses)
13   ) +
14     geom_bar()

```



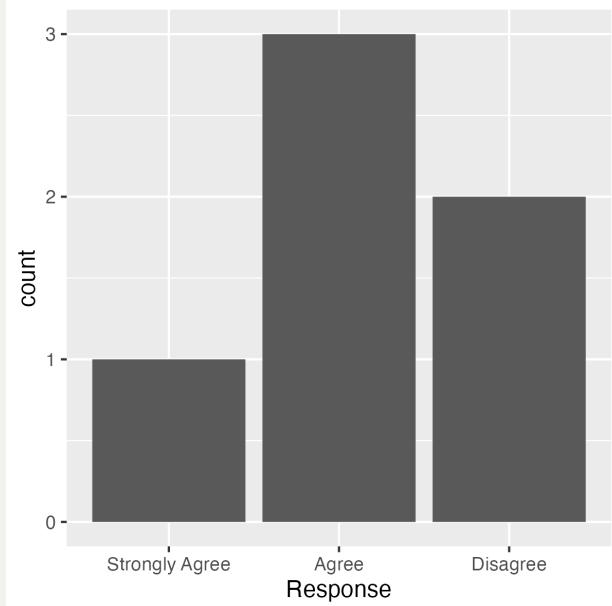
By default, `factor` is sorted alphabetically.

`ggplot2` also plots `factor` in that order, so you can adjust the `factor` to adjust the drawing order.

```

1 res =
2   data.frame(responses)
3 # Sort by level of agreement from strong -> weak
4 res$res =
5   factor(
6     res$res,
7     levels =
8       c("Strongly Agree", "Agree", "Disagree")
9   )
10
11 response_barplot2 =
12   ggplot(
13     data = res,
14     aes(x = res)
15   ) +
16   geom_bar() +
17   xlab("Response")

```



You can also use the parameter `ordered` to let others know that your `factor` is ordered properly.

```

1 responses =
2   factor(
3     c("Agree", "Agree", "Strongly Agree", "Disagree", "Disagree",
4       "Agree"),
5     ordered = TRUE
6   )

```

```

> is.ordered(responses)
[1] TRUE

```

Using `factor` to change values

You can use `recode()` in `dplyr` package to change `value`

`dplyr` is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- `mutate()` adds new variables that are functions of existing variables
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.

- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows.

These all combine naturally with `group_by()` which allows you to perform any operation “by group”. You can learn more about them in `vignette("dplyr")`. As well as these single-table verbs, dplyr also provides a variety of two-table verbs, which you can learn about in `vignette("two-table")`.

Based on the introduction on the official website of dplyr.

Here's an example:

```

1 x =
2   factor(
3     c("alpha", "beta", "gamma", "theta", "beta", "alpha")
4   )
5
6 x =
7   recode(
8     x,
9     alpha = "a",
10    beta = "b",
11    gamma = "c",
12    theta = "d"
13  )

```

The screenshot shows the RStudio interface with the code editor and the R console. The code in the editor is identical to the one above. In the R console, the user has run the code, and the output is:

```

> str(x)
Factor w/ 4 levels "a","b","c","d": 1 2 3 4 2 1
>

```

Delete useless levels

```

1 mouse.genes =
2   read.delim(
3     file = "data/talk04/mouse_genes_biomart_sep2018.txt",
4     sep = "\t",
5     header = T,
6     stringsAsFactors = T
7   )

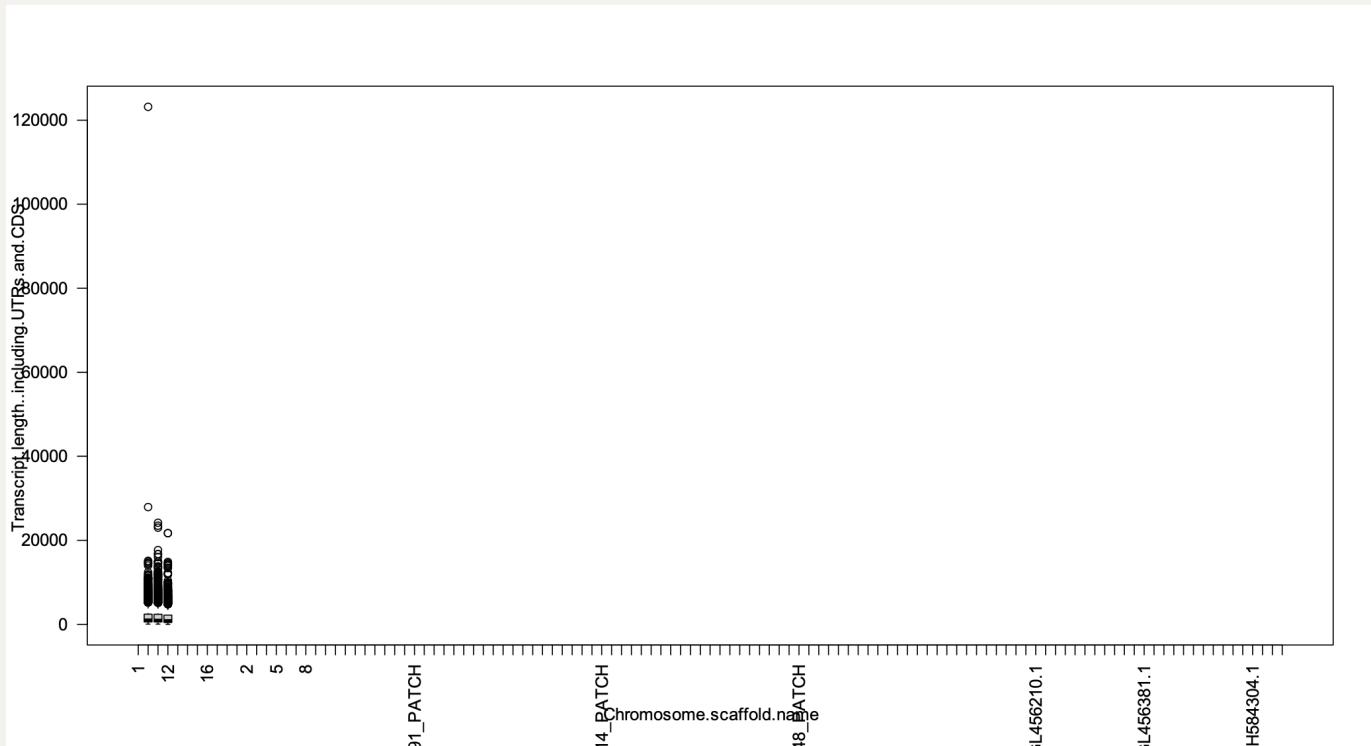
```

```

+ }
> str(mouse.genes)
'data.frame': 138532 obs. of 6 variables:
 $ Gene.stable.ID           : Factor w/ 55029 levels "ENSMUSG00000000001",...: 17960 17959 17958 17957 17956 ...
 $ Transcript.stable.ID     : Factor w/ 138532 levels "ENSMUST00000000001",...: 17312 17311 17310 17309 17308 ...
 $ Protein.stable.ID        : Factor w/ 65897 levels "", "ENSMUSP00000000001",...: 1 1 16260 1 16259 16258 1 1 1 ...
 $ Transcript.length..including.UTRs.and.CDS.: int 67 67 1144 69 519 1824 71 59 67 1378 ...
 $ Transcript.type          : Factor w/ 48 levels "3prime_overlapping_ncRNA",...: 18 18 24 18 24 24 18 18 24 ...
 $ Chromosome.scaffold.name : Factor w/ 117 levels "1", "10", "11", ...: 115 115 115 115 115 115 115 115 ...

```

If you draw a plot without deleting the useless `levels`, you will get this result:



But when you delete the useless `level` using these commands:

```

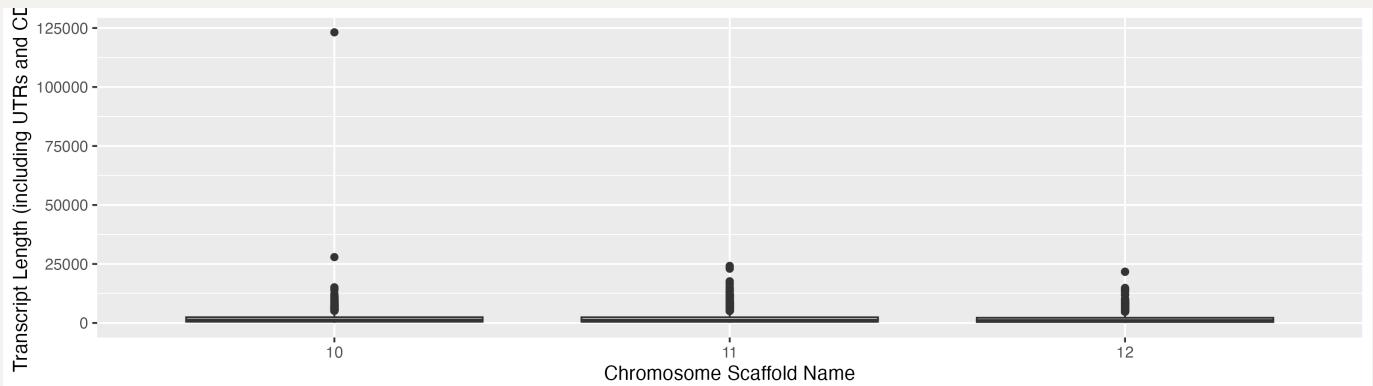
1 mouse.chr_10_12$Chromosome.scaffold.name =
2   droplevels(mouse.chr_10_12$Chromosome.scaffold.name)

```

You will see that:

```
> +  droplevels( mouse.chr_10_12$Chromosome.scaffold.name )
> levels(mouse.chr_10_12$Chromosome.scaffold.name)
[1] "10" "11" "12"
```

Then, you'll get the plot like this:



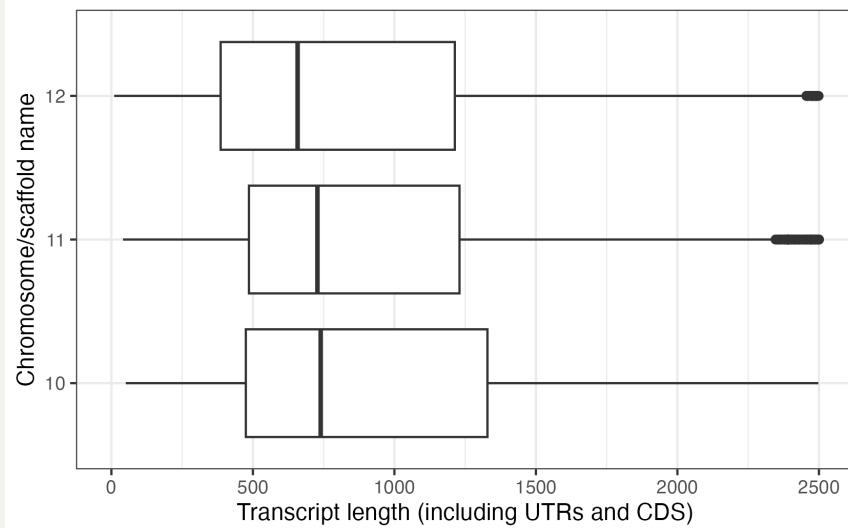
Source code:

```
1 mouse_gene_plot02 =
2   ggplot(
3     mouse.chr_10_12,
4     aes(
5       x = Chromosome.scaffold.name,
6       y = Transcript.length..including.UTRs.and.CDS.
7     )
8   ) +
9   geom_boxplot() +
10  labs(
11    x = "Chromosome Scaffold Name",
12    y = "Transcript Length (including UTRs and CDS)"
13  )
```

You can also use `tibble` to solve these problems:

```
1 mouse.tibble =
2   read_delim(
```

```
3     file = "data/talk04/mouse_genes_biomart_sep2018.txt",
4     delim = "\t",
5     quote = "",
6     show_col_types = FALSE
7 )
8
9 mouse.tibble.chr10_12 =
10   mouse.tibble %>% filter(
11     `Chromosome/scaffold name` %in% c("10", "11", "12"))
12
13
14 mouse_gene_plot03 =
15   ggplot(
16     mouse.tibble.chr10_12,
17     aes(
18       x = Chromosome.scaffold.name,
19       y = Transcript.length..including.UTRs.and.CDS.
20     )
21   ) +
22   geom_boxplot() +
23   labs(
24     x = "Chromosome",
25     y = "Transcript length (bp)"
26   ) +
27   coord_flip() +
28   ylim(0, 2500) +
29   theme_bw()
```



Advance usage

- Use `reorder()` function to reorder the level.

```

1 x = reorder(
2   `Chromosome/scaffold name`,
3   `Transcript length (including UTRs and CDS)` ,
4   median
5 )

```

- Use `forcats::fct_reorder()` to reorder factors

```

1 x = fct_reorder(
2   `Chromosome/scaffold name`,
3   `Transcript length (including UTRs and CDS)` ,
4   median
5 )

```