

Hertie Coding Club



Session 2: (Re) Introduction to R

Jorge Roa

Agenda for today

Agenda for today

- Little introduction
- Recap from the first session
- Tidyverse and packages
- Base R and Tidyverse
- Wrangle dataframes

Little Introduction

Little Introduction

I want to know you:

I want you to answer four questions:

- Your name
- What are you studying?
- What would you sing at Karaoke night?
- Any unpopular opinion?



Recap from the first session

Recap from the first session

Objects:

To create an object, give it a name followed by the assignment operator, followed by the value.

- Assignment operator `<--`
- Can also use `=` but not recommended
- Shortcut: `Alt + -` on PC, `Option + -` on Mac

```
1 x <- 2 + 2  
2  
3 x
```

```
[1] 4
```


More type of objects

There are 5 basic types of objects in the R language:

- **Atomic vectors** are one of the basic types of objects in R programming. Atomic vectors can store homogeneous data types such as character, doubles, integers, raw, logical, and complex.
- **List** is another type of object in R programming. List can contain heterogeneous data types such as vectors or another lists.

```
1 #Numeric vector
2 numbers <- c(1, 2, 3, 4)
3
4 #String vector
5 characters <- c("a", "b", "c", "d")
6
7 #Numeric value
8 value <- 5
9
10 #List
11 my_list <- list(c(1, 2, 3, 4), list("a", "b", "c"))
```

```
1 print(numbers)
```

```
[1] 1 2 3 4
```

```
1 print(characters)
```

```
[1] "a" "b" "c" "d"
```

```
1 print(value)
```

```
[1] 5
```

```
1 print(my_list)
```

```
[[1]]
```

```
[1] 1 2 3 4
```

```
[[2]]
```

```
[[2]][[1]]
```

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

```
[[2]][[2]]
```

```
[1] "b"
```

```
[[2]][[3]]
```

```
[1] "c"
```


More type of objects

- **Matrices:** To store values as 2-Dimensional array, matrices are used in R. Data, number of rows and columns are defined in the `matrix()` function.
- **Factors:** Factor object encodes a vector of unique elements (levels) from the given data vector.
- **Arrays:** `array()` function is used to create n-dimensional array. This function takes `dim` attribute as an argument and creates required length of each dimension as specified in the attribute.

```
1 x <- c(1, 2, 3, 4, 5, 6)
2
3 # Matrix
4 mat <- matrix(x, nrow = 2)
5
6 # array
7 arr <- array(c(1, 2), dim = c(3, 3))
```

```
1 print(mat)
```

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

```
1 print(arr)
```

	[,1]	[,2]	[,3]
[1,]	1	2	1
[2,]	2	1	2
[3,]	1	2	1

Finally: dataframes

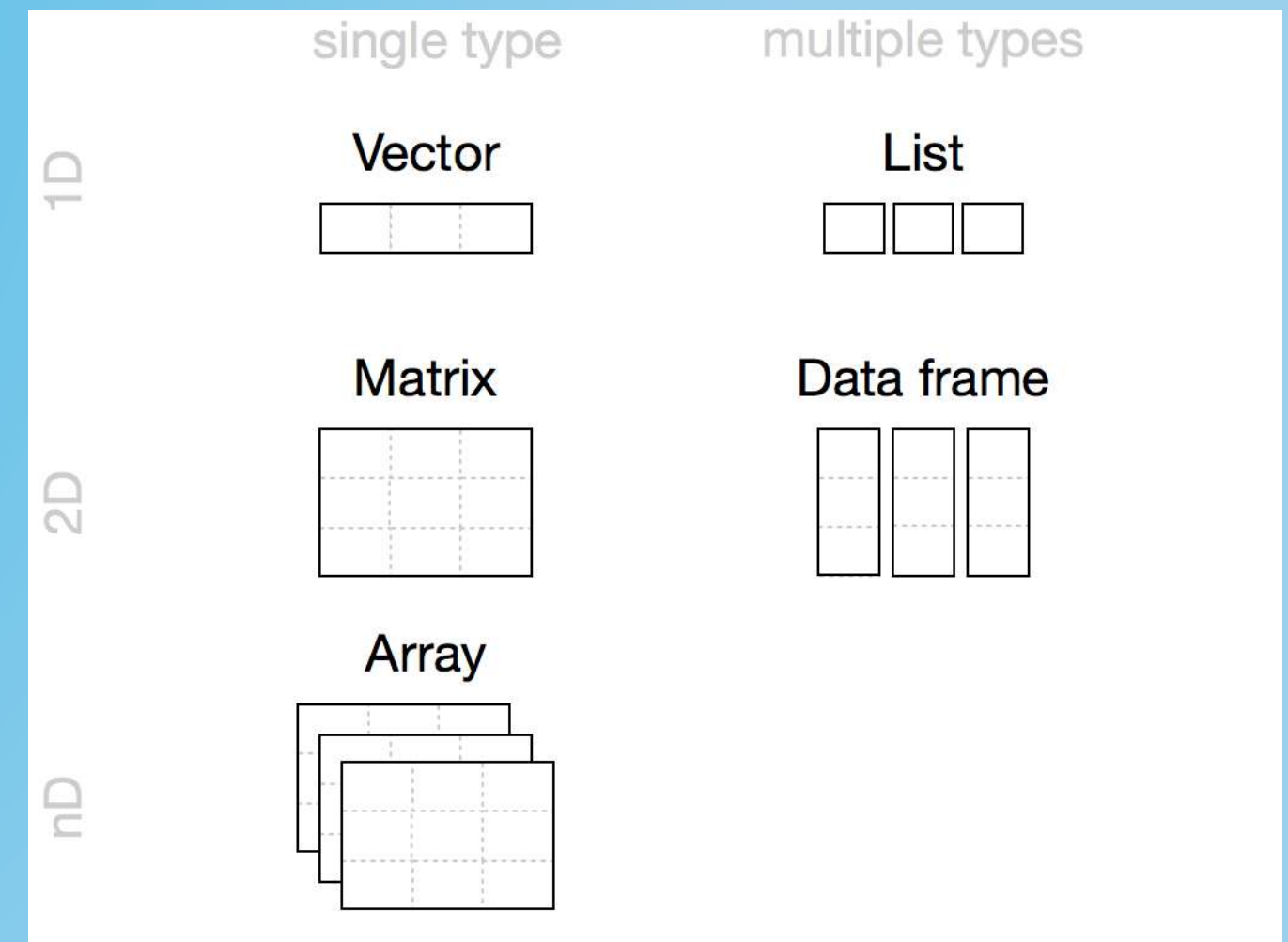
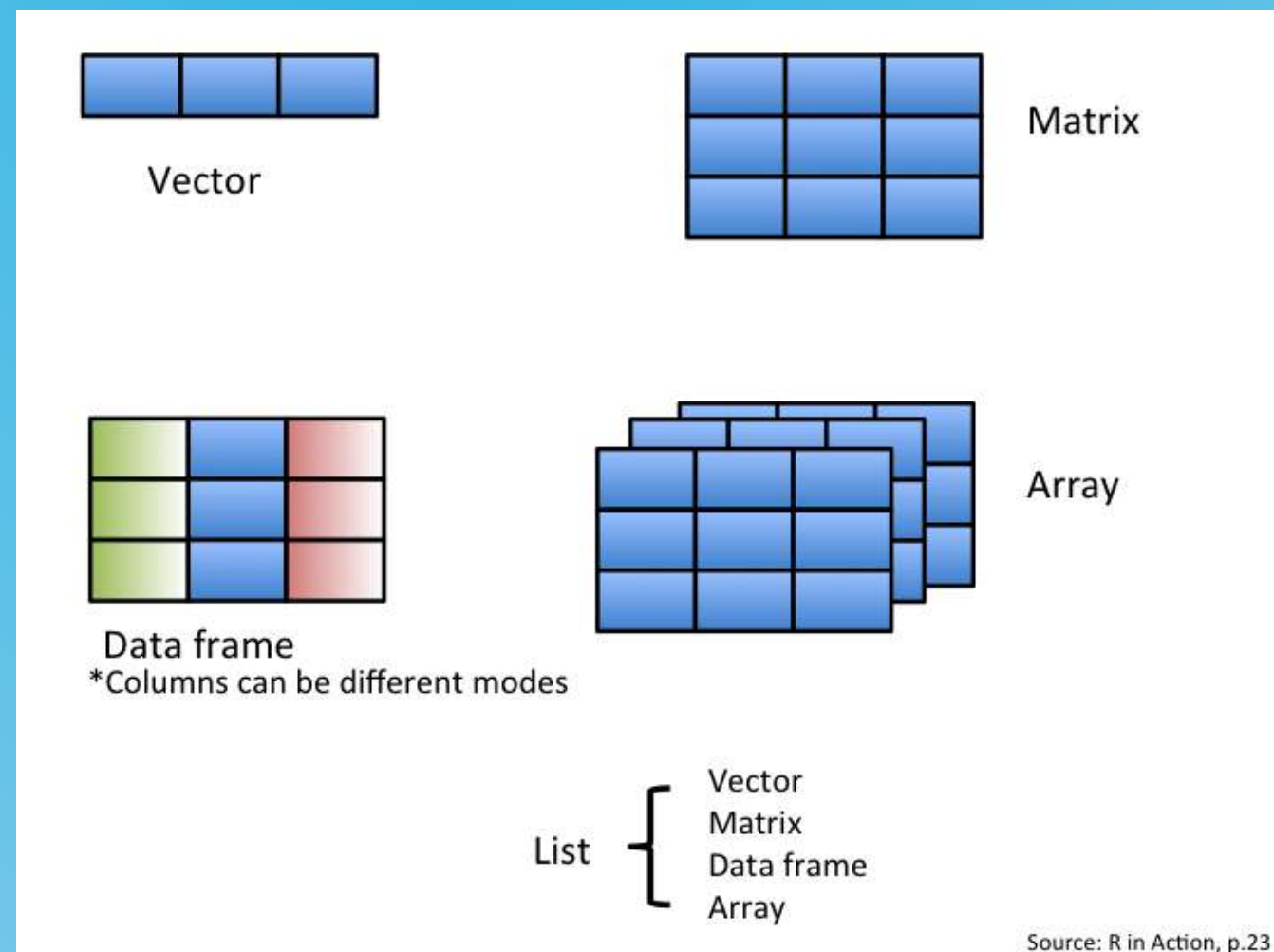
- Data frames are 2-dimensional tabular data object in R programming.
- Data frames consists of multiple columns and each column represents a vector.
- Columns in data frame can have different modes of data unlike matrices.

```
1 # Create vectors
2 who <- c("Mom", "Sister", "Myself", "Dad", "Brother", "Brother", "Our dog (:")
3 age <- c(58, 17, 25, 60, 29, 27, 5)
4 names <- c("Carmen", "Fernanda", "Jorge", "Arturo", "Ale", "Eduardo", "Rocky")
5
6 # Create data frame of vectors
7 df_my_family <- data.frame(who, age, names)
```

```
1 print(df_my_family)
```

	who	age	names
1	Mom	58	Carmen
2	Sister	17	Fernanda
3	Myself	25	Jorge
4	Dad	60	Arturo
5	Brother	29	Ale
6	Brother	27	Eduardo
7	Our dog (:	5	Rocky

Objects:summary



Operations of vectors

Operations of numeric vectors

- `length(x)`: how many elements you have in your vector.
- `sort(x, decreasing = F)`: sort your numerical values.
- `sum(x)`: returns the sum of your values.
- `min(x)`: minimum value of your numeric vector.
- `mean(x)`: mean of your numeric vector.
- `median(x)`: median
- `sd(x)`: Standard deviation.
- `var(x)`: variance of your numeric vector.
- `summary(x)`: summary of your numeric vector.

```
1 v_age <- c(22, 25, 36, 60, 15, 25, 20, 10)
2
3 length(v_age)
```

```
[1] 8
```

```
1 sort(v_age, decreasing = F)
```

```
[1] 10 15 20 22 25 25 36 60
```

```
1 sum(v_age)
```

```
[1] 213
```

```
1 min(v_age)
```

```
[1] 10
```

```
1 mean(v_age)
```

```
[1] 26.625
```

```
1 median(v_age)
```

```
[1] 23.5
```

```
1 sd(v_age)
```

```
[1] 15.50979
```

```
1 var(v_age)
```

```
[1] 240.5536
```

```
1 summary(v_age)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
10.00	18.75	23.50	26.62	27.75	60.00

Exercises

Exercises

```

1 #1.-Create a list containing strings, numbers, vectors and a logical values.
2
3 #2.-Create a dataframe of 5 variables
4 #Hint: Remember the length of the vectors
5
6 #3.- Create a vector with numerical values and strings with a length of 10
7
8 #4.- Assign the following vectors to a meaningful variable name:
9 #Hint: Remember the assignment operator.
10
11 c(2, 4, 6, 8, 10, 12, 14, 16, 20)
12 0
13 3.141593
14 c(1, 10, 100, 1000, 10000, 100000)
15
16 #5.- Create vectors that correspond to the following variables names:
17
18 yourage
19 days_of_the_week
20 firstFivePrimeNumbers

```

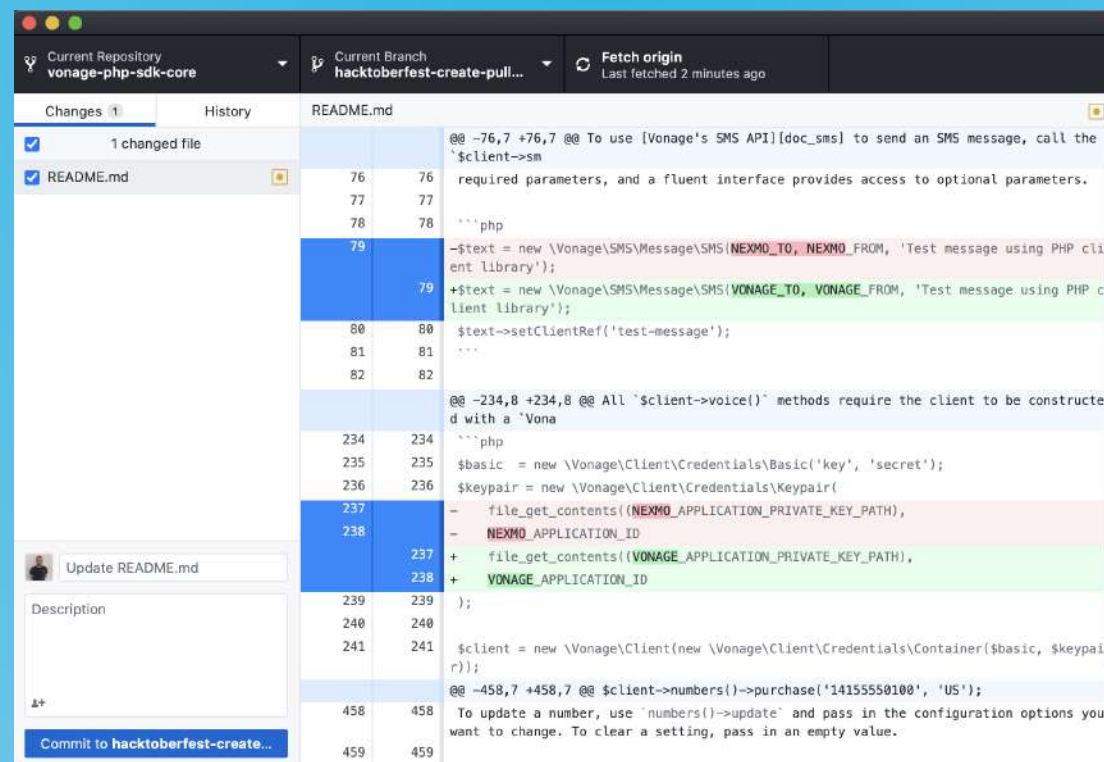
```

1 #1.-Create a list containing strings, numbers, vectors and a logical values.
2 list <- list(c("Coding", "Club"), c(1,2,3), 7)
3
4 #2.-Create a dataframe of 5 variables
5 #Hint: Remember the length of the vectors
6 df_my_family <- data.frame(number = c(1,2,3),
7                             age = c(17,18,19),
8                             name = c("Alex", "Eduardo", "Jorge"),
9                             favorite_color = c("blue", "orange", "black"),
10                             favorite_number = c(20, "5", 50))
11
12 #3.- Create a vector with numerical values and strings with a length of 10
13 vector <- c(1,2,3,"number",99, 100, "yes", "hi", 9, 10)
14
15 length(vector)
16
17 #4.- Assign the following vectors to a meaningful variable name:
18 #Hint: Remember the assignment operator.
19 vector <- c(2, 4, 6, 8, 10, 12, 14, 16, 20)
20 value <- 0
21 value_2 <- 3.141593
22 num_vector <- c(1, 10, 100, 1000, 10000, 100000)
23
24 #5.- Create vectors that correspond to the following variables names:
25 yourage <- c(25)
26 days_of_the_week <- c("Monday", "Tuesday", "Wednesday",
27                       "Thursday", "Friday", "Saturday",
28                       "Sunday")
29 firstFivePrimeNumbers <- c(2,3,5,7,11)

```


GitHub

GitHub

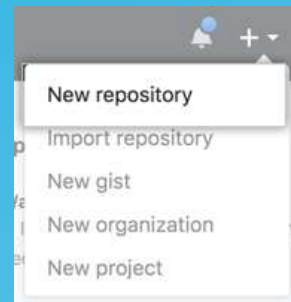


- Website and cloud-based service to store and manage code
- Git IDE: used in the programming world. It is used for tracking changes in the source code during software development.
- It makes it easier for individuals and teams to use Git for version control and collaboration.

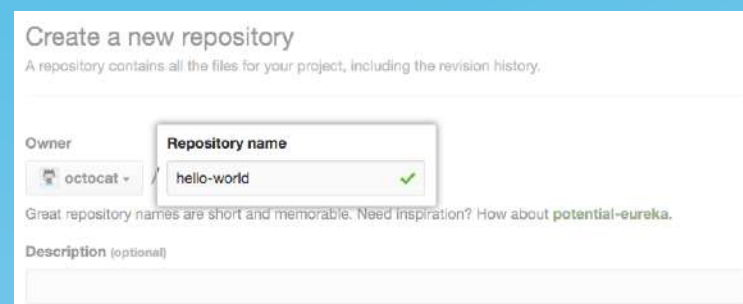
Create our first repo

Create our first repo

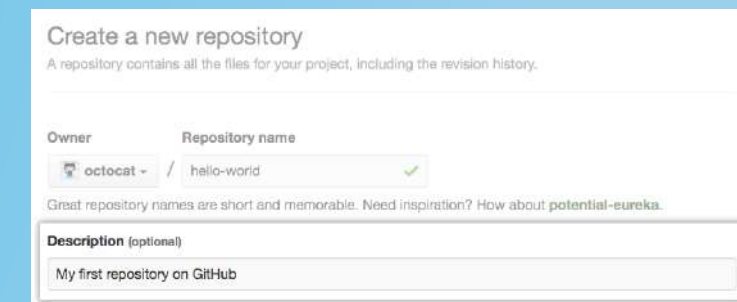
1. In the upper-right corner of any page, use the drop-down menu, and select New repository.



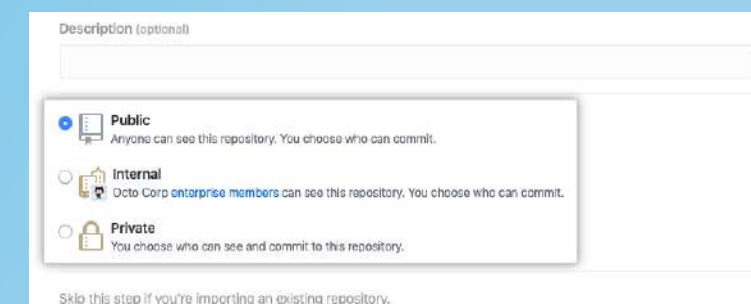
2. Type a short, memorable name for your repository. For example, “intro-to-r”.



3. Optionally, add a description of your repository. For example, “My first repository on GitHub.”

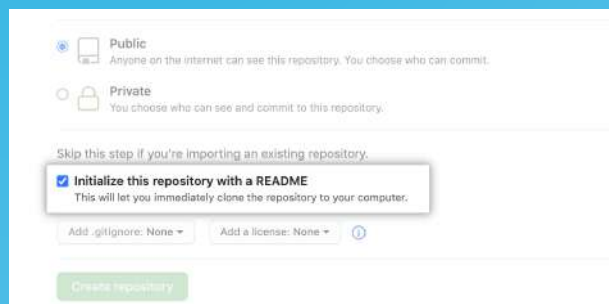


4. Choose a repository visibility.



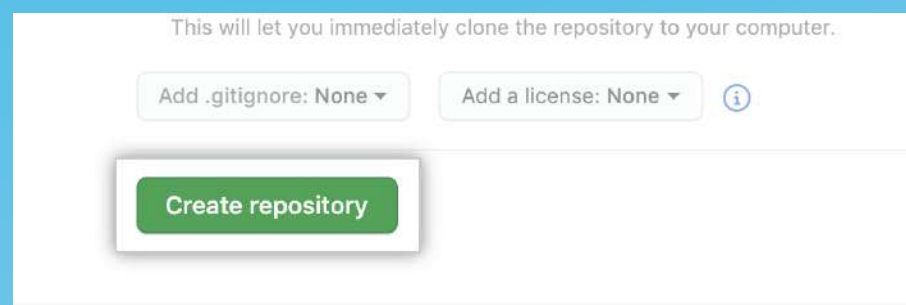
Create our first repo

5. Select Initialize this repository with a README.



The screenshot shows the GitHub repository creation interface. It has two radio buttons for visibility: 'Public' (selected) and 'Private'. Below them is a section for initialization with a README, which is checked. There are also dropdowns for adding a .gitignore file and a license, both set to 'None'. A green 'Create repository' button is at the bottom.

6. Click Create repository.



This screenshot is a close-up of the bottom part of the GitHub repository creation form. It shows the text 'This will let you immediately clone the repository to your computer.' above two dropdown menus: 'Add .gitignore: None' and 'Add a license: None'. Below these is a prominent green 'Create repository' button.



CONGRATULATIONS: YOU CREATED YOUR FIRST REPO

We will explain how it works

Tidyverse and packages

Tidyverse and packages

- R packages are a collection of R functions, complied code and sample data.
- **Why:** There are millions of functions. If they were all preloaded, there wouldn't be enough RAM to work with. There are packages of such varied disciplines that we likely use relatively few.
- They are stored under a directory called “library” in the R environment.
- By default, R installs a set of packages during installation. More packages are added later, when they are needed for some specific purpose.
- When we start the R console, only the default packages are available by default.
- Other packages which are already installed have to be loaded explicitly to be used by the R program that is going to use them.
- We can also generate our functions and even create an R package!



Tidyverse

Package set for: Import, Clean, Transform, Process, Analyze and Visualize



File Import/Export

readr

Package set for: load plain text files (txt, csv, tsv)



File Import/Export

readxl

Package set for: load excel files (xls, xlsx)



File Import/Export

haven

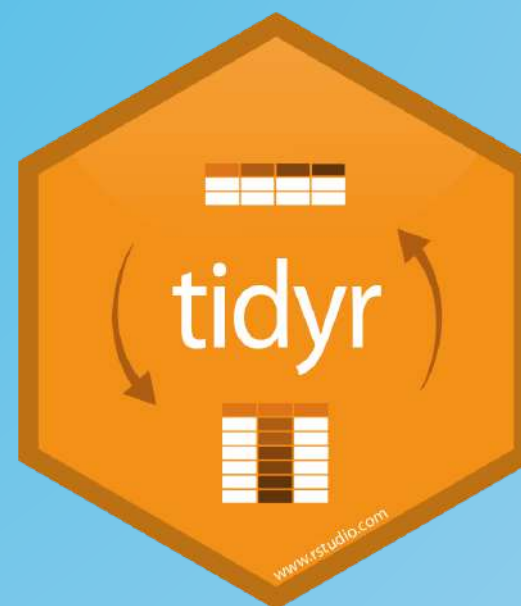
Package set for: Display proprietary formats (dta, sav). Like STATA and other formats.



Wrangling data

tidyr

Package set for: transform dataframe structures



Wrangling data

lubridate

Package set for: wrangling dates. Tools that make working with dates and times easier.



Wrangling data

stringr

Package set for: wrangling string or characters.



Wrangling data

dplyr

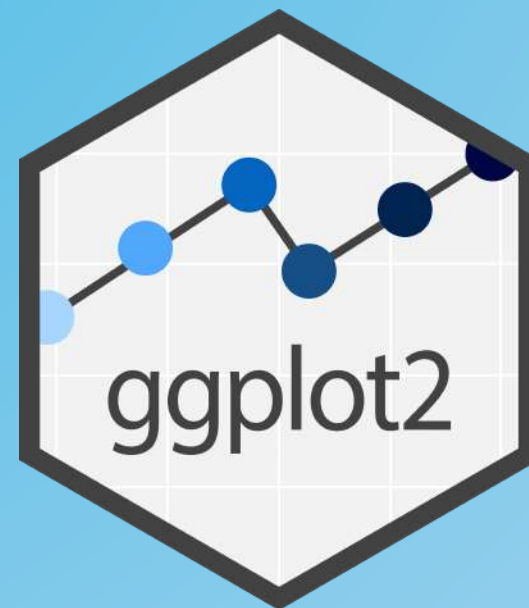
Package set for: wrangling dataframes. facilitates several functions for the data frames in R. dplyr package is for data wrangling and data analysis purposes.



Analysis and visualization

ggplot

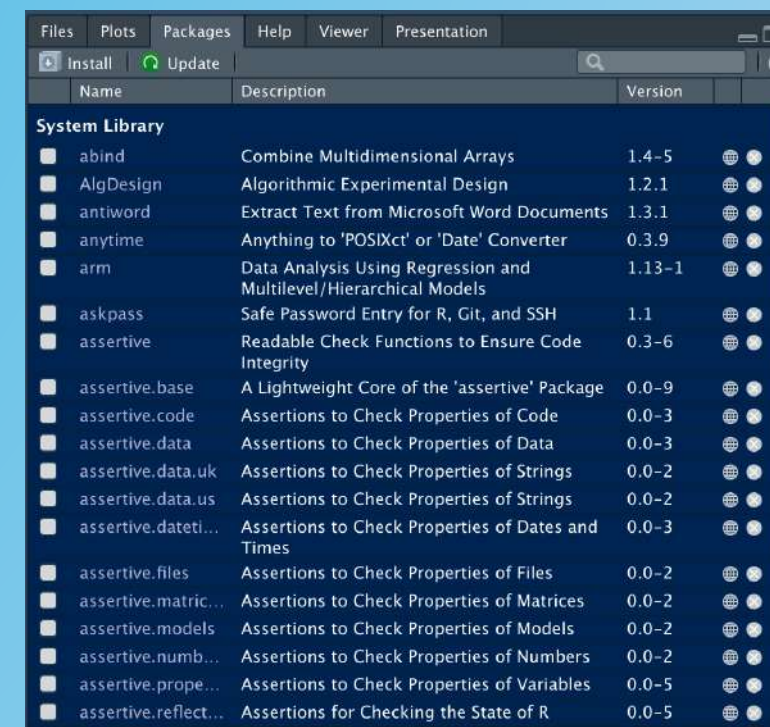
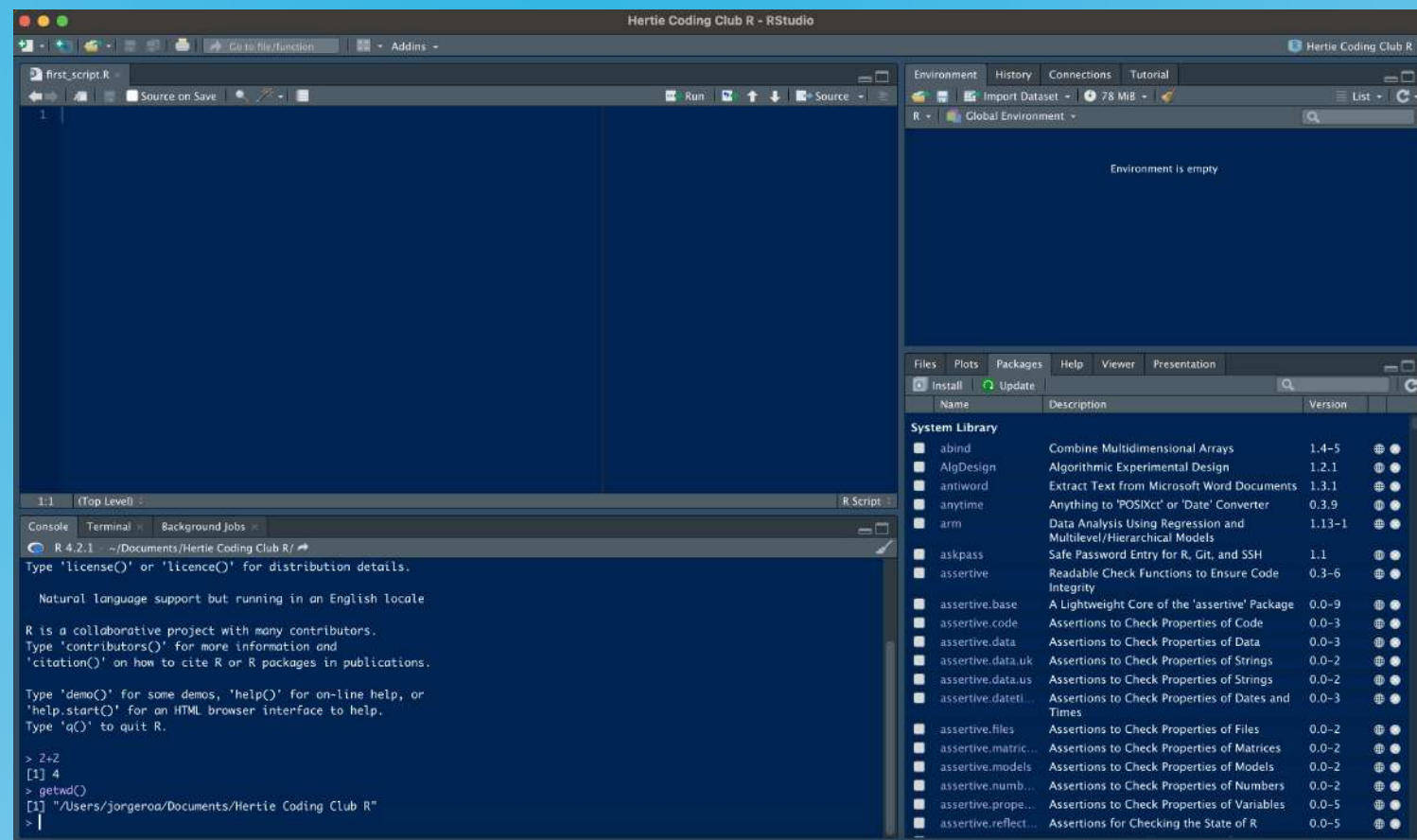
Package set for: plots and maps. One of the most popular visualization package in R.



How we install packages

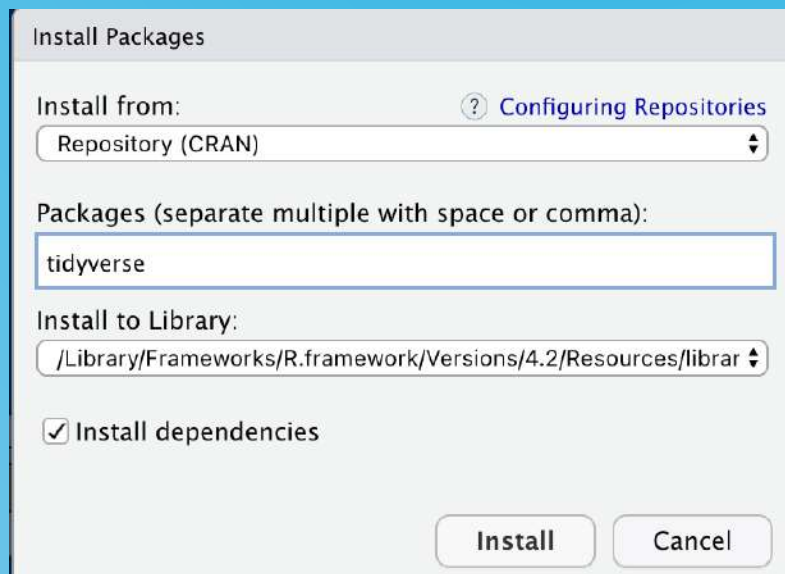
The easy way

- Go to the “Packages” tab
- Press the “Install” button



How we install packages

The easy way



Install Packages

Install from: [? Configuring Repositories](#)

Repository (CRAN)

Packages (separate multiple with space or comma):

tidyverse

Install to Library:

/Library/Frameworks/R.framework/Versions/4.2/Resources/librar

☒ Install dependencies

Install Cancel

- Other way is type in the console
- `install.packages("tidyverse")`

```
> install.packages("tidyverse")
trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.2/tidyverse_1.3.2.tgz'
Content type 'application/x-gzip' length 420896 bytes (411 KB)
=====
downloaded 411 KB

The downloaded binary packages are in
  /var/folders/q3/ztffc4r10tv5h4pftjr2qw0r0000gn/T//RtmpKwSpEV/downloaded_packages
>
```

Dataframes

Dataframes

Example

- We will work with Airbnb accommodation data in Berlin as of September 15, 2022. They are open data available at [Airbnb: get the data](#).
- They are open data licensed under the Creative Commons CC0 1.0 Universal “Public Domain Dedication.”
- Those who stay can choose between entire houses/apartments, only private rooms, or shared rooms (room_type).
- After the stay, they must leave an evaluation (review).
- Accommodations vary in price, a minimum number of days of stay, days available, etc.



Import dataframes from your computer

How we import data?



Excel files

Quite frequently, the sample data is in Excel format, and needs to be imported into R prior to use.

```
1 library(readxl)
2
3 df_listings <- read_xlsx("data/listings.xlsx")
```

id <dbl>	listing_url <chr>	scrape_id <dbl>	last_scraped <dtm>
6.528688e+17	https://www.airbnb.com/rooms/652868795892201022	2.022092e+13	2022-09-15
2.907769e+07	https://www.airbnb.com/rooms/29077694	2.022092e+13	2022-09-16
2.708061e+07	https://www.airbnb.com/rooms/27080612	2.022092e+13	2022-09-15
6.656484e+17	https://www.airbnb.com/rooms/665648367391379321	2.022092e+13	2022-09-15
3.176000e+03	https://www.airbnb.com/rooms/3176	2.022092e+13	2022-09-15
3.773800e+07	https://www.airbnb.com/rooms/37738004	2.022092e+13	2022-09-16
3.309000e+03	https://www.airbnb.com/rooms/3309	2.022092e+13	2022-09-15
7.071000e+03	https://www.airbnb.com/rooms/7071	2.022092e+13	2022-09-16
2.019132e+06	https://www.airbnb.com/rooms/2019132	2.022092e+13	2022-09-16
9.991000e+03	https://www.airbnb.com/rooms/9991	2.022092e+13	2022-09-16

1-10 of 10,000 rows | 1-4 of 75 columns

Previous 1 2 3 4 5 6 ... 1000 Next

What do we want to find out from this data?

First, I ask myself questions, then think about the code that answers them.

- What are the variables? How many?
- How many observations do you have?
- What values do these variables take?
- Is there missing data? Are there duplicate cases?

Explore the data

```
1 View(df_listings)
```

OR (Best option: only works with dataframes)

CTRL + click on your object (IN YOUR SCRIPT).

	id	listing_url	scrape_id	last_scraped	source	name
1	6.528688e+17	https://www.airbnb.com/rooms/6528687958922010...	2.022092e+13	2022-09-15	city scrape	Kleine Auszeit? Oder Business-Trip? Alles möglich!
2	2.907769e+07	https://www.airbnb.com/rooms/29077694	2.022092e+13	2022-09-16	previous scrape	Wohnung im Gr- ⁿ en nah an der Metropole
3	2.708061e+07	https://www.airbnb.com/rooms/27080612	2.022092e+13	2022-09-15	city scrape	Apartment with Living/Sleeping Room & own Kitchen
4	6.656484e+17	https://www.airbnb.com/rooms/6656483673913793...	2.022092e+13	2022-09-15	city scrape	Sch- ⁿ e Eigentumswohnung mit Balkon
5	3.176000e+03	https://www.airbnb.com/rooms/3176	2.022092e+13	2022-09-15	city scrape	Fabulous Flat in great Location
6	3.773800e+07	https://www.airbnb.com/rooms/37738004	2.022092e+13	2022-09-16	previous scrape	Remise Villa Erica Superior Apartment
7	3.309000e+03	https://www.airbnb.com/rooms/3309	2.022092e+13	2022-09-15	city scrape	BerlinSpot Sch- ⁿ eberg near KaDeWe
8	7.071000e+03	https://www.airbnb.com/rooms/7071	2.022092e+13	2022-09-16	previous scrape	BrightRoom with sunny greenview!
9	2.019132e+06	https://www.airbnb.com/rooms/2019132	2.022092e+13	2022-09-16	previous scrape	Villa am Berliner Stadtrand f- ^r bis zu 24 G- ^u ste
10	9.991000e+03	https://www.airbnb.com/rooms/9991	2.022092e+13	2022-09-16	city scrape	Georgeous flat - outstanding views
11	1.432500e+04	https://www.airbnb.com/rooms/14325	2.022092e+13	2022-09-15	city scrape	Studio Apartment in Prenzlauer Berg
12	3.257390e+07	https://www.airbnb.com/rooms/32573899	2.022092e+13	2022-09-16	previous scrape	Modern House in Berlin/Teltow-Seehof
13	1.664400e+04	https://www.airbnb.com/rooms/16644	2.022092e+13	2022-09-16	city scrape	In the Heart of Berlin - Kreuzberg
14	1.791020e+05	https://www.airbnb.com/rooms/179102	2.022092e+13	2022-09-15	city scrape	The Special Place
15	1.804400e+05	https://www.airbnb.com/rooms/180440	2.022092e+13	2022-09-16	previous scrape	CITY STUDIO WEST@ Kurf- ^v rstendamm(monthly disc...
16	1.811600e+05	https://www.airbnb.com/rooms/181160	2.022092e+13	2022-09-15	city scrape	Berliner Flair
17	1.790400e+04	https://www.airbnb.com/rooms/17904	2.022092e+13	2022-09-16	city scrape	Beautiful Kreuzberg studio - 3 months minimum
18	1.828160e+05	https://www.airbnb.com/rooms/182816	2.022092e+13	2022-09-15	city scrape	Cozy 90sqm in Prenzlauer Berg
19	2.085800e+04	https://www.airbnb.com/rooms/20858	2.022092e+13	2022-09-16	city scrape	Designer Loft in Berlin Mitte
20	1.839880e+05	https://www.airbnb.com/rooms/183988	2.022092e+13	2022-09-15	city scrape	Luxury Ku'damm Studio + Bikes 4you

Explore the data

The `dim()`, `names()`, and `str()` functions take a data frame as an argument.

```
1 dim(df_listings) # dimension of the dataframe
[1] 16680      75

1 nrow(df_listings) # rows
[1] 16680

1 ncol(df_listings) # columns (variables)
[1] 75

1 names(df_listings)[1:10] # variable names
[1] "id"           "listing_url"   "scrape_id"
[4] "last_scraped" "source"        "name"
[7] "description"  "neighborhood_overview" "picture_url"
[10] "host_id"

1 #str(df_listings)
```

Explore the data

```
1 max(df_listings$price) #maximum value.
[1] 4375

1 median(df_listings$price) #median.
[1] 65

1 min(df_listings$price) #minimum value.
[1] 0

1 mean(df_listings$price) #mean.
[1] 96.30809

1 median(df_listings$price) #median.
[1] 65

1 var(df_listings$price) #variance.
[1] 13638.87

1 sd(df_listings$price) #Standard deviation.
[1] 116.7856
```

Exercise

- Upload the Barcelona listings.xlsx file.
- Calculate the mean, median, and variance of the variables `minimum_nights`, `number_of_reviews`, and `last_review`. Then, store those results in different vectors and create a dataframe with them. Name the dataframe with `df_exercise_1`

Exercise

- Upload the Berlin listings.csv file.
- Calculate the mean, median, and variance of the variables `minimum_nights`, `number_of_reviews`, and `last_review`. Then, store those results in different vectors and create a dataframe with them. Name the dataframe with `df_exercise_1`
- Why get NA in the variance of the `last_review` variable?
 - It's a string. We can't apply numerical function to strings.

```

1 v_mean_minimum_nights <- mean(df_listings$minimum_nights)
2 v_mean_number_of_reviews <- mean(df_listings$number_of_reviews)
3 v_mean_last_review <- mean(df_listings$last_review)
4
5 v_med_minimum_nights <- median(df_listings$minimum_nights)
6 v_med_number_of_reviews <- median(df_listings$number_of_reviews)
7 v_med_last_review <- median(df_listings$last_review)
8
9 v_var_minimum_nights <- var(df_listings$minimum_nights)
10 v_var_number_of_reviews <- var(df_listings$number_of_reviews)
11 v_var_last_review <- var(df_listings$last_review)
12
13
14 df_exercise_1 <- data.frame(mean = c(v_mean_minimum_nights,
15                                   v_mean_number_of_reviews,
16                                   v_mean_last_review),
17                             median = c(v_med_minimum_nights,
18                                       v_med_number_of_reviews,
19                                       v_med_last_review),
20                             variance = c(v_var_minimum_nights,
21                                         v_var_number_of_reviews,
22                                         v_var_last_review))
23
24
25 df_exercise_1

```

	mean	median	variance
1	12.11559	3	1530.456
2	27.76787	6	3881.410
3	NA	NA	NA

Thanks for your time



Remember that everybody can learn how to code!!

1,100 lines of code were created for this presentation.