# Introduction to R

## Session 02: Working with data in R

Álvaro Pérez[1]

Instituto Tecnológico Autónomo de México

Fall 2024

---

[1]Based on PhD Romero Londoño's notes.

# Working with Data

▶ R is all about working with data!

▶ data.frames:

   ▶ data.frames are an object type

   ▶ Most of the time, you'll be doing calculations using them

   ▶ Conceptually, data.frames are basically spreadsheets

   ▶ Technically, they're a list of vectors

ITAM

▶ It's a collection of vectors of the same length

▶ Note the use of $=$ here, not $<$-

```
1  df <- data.frame(
2  RacePosition = 1:5,
3  WayTheySayHi = as.factor(c('Hi','Hello','Hey','Yo','Hi
       ')),
4  NumberofKids = c(3,5,1,0,2))
5  df
```

▶ If we just want a quick overview:

    ▶ Arrow in the Environment tab

    ▶ 'head()' (look at the head of the data - first six rows)

    ▶ 'str()' (structure)

```
1   str(df)
```

▶ Now we have a data frame, 'df'. How do we use it?

▶ We can pull the vectors back out with '$' Note
autocompletion of variable names.

▶ We can treat it just like the vectors we had before

```
1  df$NumberofKids
2  df$NumberofKids[2]
3  df$NumberofKids >= 3
```

# Practice

▶ Create 'df2 <- data.frame(a = 1:20, b = 0:19*2,' 'c = sample(101:200,20,replace=TRUE))'

▶ What is the average of 'c'?

▶ What is the sum of 'a' times 'b'?

▶ Did you get any values of 'c' 103 or below? (make a logical)

▶ What is on the 8th row of 'b'?

▶ How many rows have 'b' above 10 AND 'c' below 150?

# Practice Answers

```
1  df2 <- data.frame(
2  a = 1:20,
3  b = 0:19*2,
4  c = sample(101:200,20,replace=TRUE))
5
6  mean(df2$c)
7
8  sum(df2$a*df2$b)
9
10 sum(df2$c <= 103) > 0
11
12 df2$b[8]
13
14 sum(df2$b > 10 & df2$c < 150)
```

Álvaro PL

- ▶ We can manipulate data frames:

    - ▶ Create new variables

    - ▶ Change variables

    - ▶ Rename variables

- ▶ It's very common that you'll have to work with data before analyzing it

- ▶ "data cleaning" is very important and a big part of statistical analysis

# Creating New Variables

▶ data.frames are just lists of vectors

▶ So create a vector and tell R where in that list to stick it

▶ Use descriptive names so you know what the variable is

```
1  df$State <- c('Alaska','California',
2                      'California','Maine',
3                      'Florida')
4  df
```

# Another approach - DPLYR and Tidyverse

- ▶ We just saw the base-R way to do it

- ▶ Can use **dplyr** (data pliers) for data manipulation instead

- ▶ dplyr syntax is inspired by SQL

- ▶ tidyverse isn't a part of base R. It's in a package, so we'll need to install it

# Variable creation with dplyr

- The **mutate** command will "mutate" our data frame to have a new column in it
- The pipe '%>%' says "take df and send it to that mutate command to use"
- Or we can stick the data frame itself in the 'mutate' command
- Thus these two are equivalent:

```
1  library(tidyverse)
2  df1 <- df %>% mutate(State = c('Alaska','California',
3                                 'California','Maine','
                                     Florida'))
4
5  df2 <- mutate(df,State = c('Alaska','California',
6                             'California','Maine','
                                 Florida'))
7
8  identical(df1,df2)
9  df <- df1
```

# Creating New Variables

▶ We can create multiple new variables in one mutate command

```
1  df <- df %>% mutate(
2  MoreThanTwoKids = NumberofKids > 2,
3  One = 1,
4  KidsPlusPosition = NumberofKids + RacePosition)
5
6  df
```

Álvaro PL

# Manipulating Variables

▶ We can't really **change** variables, but we can overwrite them

▶ We can drop variables with '-' in the dplyr 'select' command

▶ Note we chain multiple dplyr commands with '%>%'

```
1  df <- df %>%
2    select(-KidsPlusPosition,-WayTheySayHi,-One) %>%
3    mutate(State = as.factor(State),
4           RacePosition = RacePosition - 1)
5  df$State[3] <- 'Alaska'
6  str(df)
```

# Renaming Variables

- ▶ Sometimes it will make sense to change the names of the variables we have.

- ▶ Names are stored in 'names(df)' which we can edit directly

- ▶ Or the 'rename()' command in dplyr has us covered

```
1  names(df)
2  #two ways of renaming
3  #names(df) <- c('Pos','Num.Kids','State','mt2Kids
       ')
4
5  df <- df %>% rename(Pos = RacePosition, Num.Kids=
       NumberofKids,
6                       mt2Kids = MoreThanTwoKids)
7
8  names(df)
```

# Practice

▶ Create a data set 'data' with three variables: 'a' is all even numbers from 2 to 20, 'b' is 'c(0,1)' over and over, and 'c' is any ten-element numeric vector of your choice.

▶ Rename them to 'EvenNumbers', 'Treatment', 'Outcome'.

▶ Add a logical variable called Big that's true whenever EvenNumbers is greater than 15

▶ Increase Outcome by 1 for all the rows where Treatment is 1.

▶ Create a logical AboveMean that is true whenever Outcome is above the mean of Outcome.

▶ Display the data structure

# Practice Answers

```
1  data <- data.frame(a = 1:10*2,
2                     b = c(0,1),
3                     c = sample(1:100,10,replace=FALSE))
                        %>%
4    rename(EvenNumbers = a, Treatment = b, Outcome = c)
5
6  data <- data %>%
7    mutate(Big = EvenNumbers > 15,
8           Outcome = Outcome + Treatment,
9           AboveMean = Outcome > mean(Outcome))
10
11 str(data)
```

# dplyr package

- ▶ select(): used to select columns of a data frame that you want to focus on
- ▶ filter(): used to extract subsets of rows from a data frame
- ▶ arrange(): used to reorder rows of a data frame according to one of the variables/columns
- ▶ rename(): renaming a variable
- ▶ mutate(): compute transformations of variables in a data frame
- ▶ group_by():used to generate summary statistics from the data frame within strata defined by a variable.
- ▶ %>%: stringing together multiple dplyr functions in a sequence of operations.
- ▶ summarize()

# Example

- ▶ Download the database "1976-2020-president.csv" from my github repository
- ▶ In R, upload the dataframe
- ▶ Preserve only the votes for the Democratic and the Republican Party from de 1976 elections
- ▶ Create a new dataframe with only the state, the party, the votes for each party and the total of votes
- ▶ Sort states by number of votes
- ▶ Rename de variable "party_simplified.ªs "party"
- ▶ Create new variable that is the percentage of the vote (use the dataframe of the second step)
- ▶ Group by party and calculate de total aumount of votes for each party in that election

Álvaro PL

# Answers

```
1  #install.packages("dplyr")
2  library(dplyr)
3  # Rename database
4  US_elections <- X1976_2020_president
5  head(US_elections)
6  str(US_elections)
7  # filter()
8  elect_1976 <- filter(US_elections, year == 1976 & (
       party_simplified == "DEMOCRAT" | party_simplified
       == "REPUBLICAN"))
9  # select()
10 votes_1976 <- select(elect_1976,state | party_
       simplified | candidatevotes)
11 # arrange()
12 arrange(votes_1976,candidatevotes)
13 # rename()
14 votes_1976 <- rename(votes_1976,"party"=party_
       simplified)
```

```
1
2  # mutate()
3  elect_1976 <- elect_1976 %>%
4    mutate(per_votes = candidatevotes/totalvotes)
5  # group_by()
6  votes_1976_party <- votes_1976 %>%
7    group_by(party) %>%
8    summarise(sum_votes=sum(candidatevotes,na.rm = T))
```