



C91AR | ADVANCED STATISTICS USING R

Computer Lab 4: Tidying realistic data

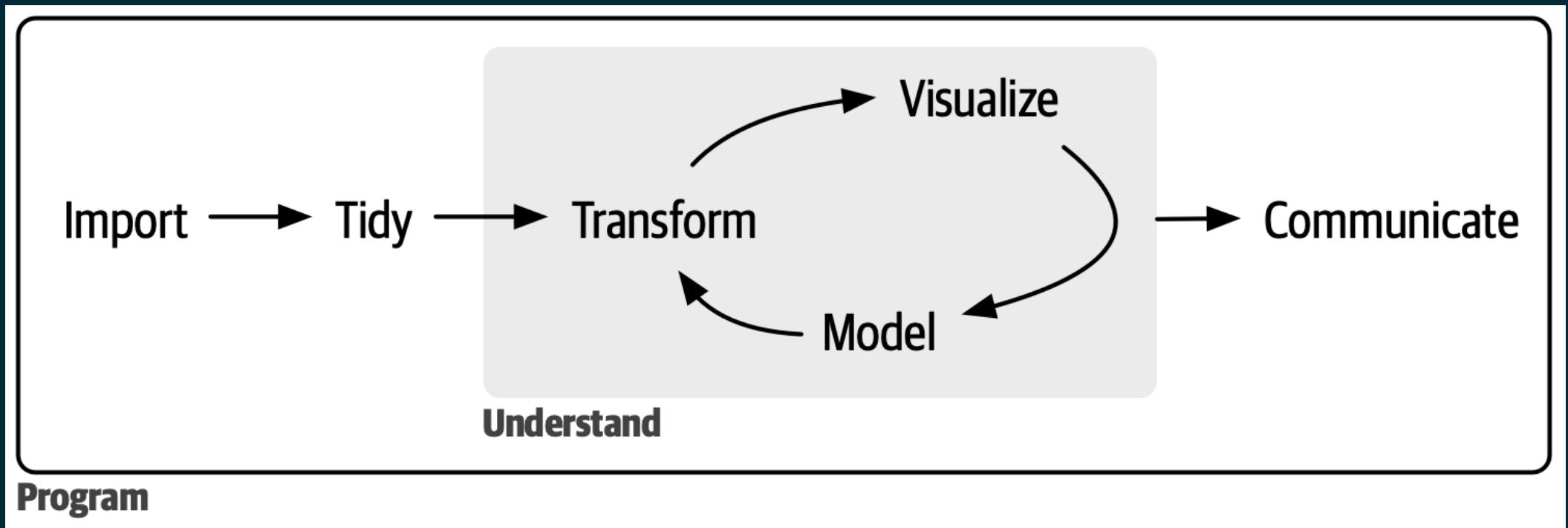
Dr Pete McKenna

2025-02-05

SETUP FOR THIS MODULE

- Continue working on R Markdown called “Lecture4_tiding_real_data” (or create this file if you weren’t here on Monday)
- Load the following packages

```
1 # load packages
2 pacman::p_load(tidyverse,
3                  snakecase,
4                  psych,
5                  summarytools)
```



SAVE THE RAW DATA PROVIDED

- Save the data from Canvas Module **Course datasets** called “robot-expression-data_raw.csv” into the “data_raw” folder

READ IN AND EXPLORE THE ROBOT EXPRESSION DATA

```
1 # Read in raw data
2 df <-
3   read_csv("data_raw/robot-expression-data_raw.csv")
```

EXPRESSION DATASET METADATA

- `ID_` = participant id
- `sex` = participant sex
- `eth` = participant ethnicity
- `eng` = “Is English your native language?”
- `Exp.date` = date of experiment
- `interact.before` = “Have you ever interacted with a robot before?
- `Alyx.is` = participants guess of robot’s gender
- `item` = plastic food item offered to robot
- `Expression` = a numeric value for each of the four robot facial expressions
- `resp` = name of the box where participants placed the food item after robot expression
- `Time` = robots onboard clock timer
- `resp_acc` = accuracy of response according to theoretically deigned expression

DATA TYPES FOR THE `col_types` ARGUMENT

- “f” = factor (i.e., categorical variable)
- “l” = logical
- “i” = integer (no decimals)
- “d” = double numeric (has decimals)
- “c” = character
- “D” = date
- “_” = skip
- “?” = guess
- For more information see the `read_csv` and `col_types` documentation in the help section
- You can pull this up by running `?read_csv` into the console
 - I’d recommend consulting the descriptor info at the top then skipping to the examples at the bottom of the help sheet

SETTING COL TYPES WITHIN THE `read_csv` FUNCTION

```
1 # Read in data and supply data type argument
2 df <-
3   read_csv("data_raw/robot-expression-data_raw.csv",
4           col_types = "ffff?ffffftf") # this is good for Base R
5
6 df <-
7   read_csv("data_raw/robot-expression-data_raw.csv",
8           col_types = "cccc?ccccctc") # tidyverse prefers character vector
```

- Now the data from the csv is stored in an object called `df`

DATA EXPLORATION

- Here are a few handy functions to exploring your data
 - `summary` : from base R
 - `glimpse` : from *tidyverse*
 - `headTail` : from *psych*
 - `dfSummary`: from *summarytools*

SUMMARISING DATA FUNCTION

- The `summary` function summarises continuous and categorical vectors
- NOTE: `summary` will not count character vectors, that's why we specified that our categorical vectors were factor type (e.g., "f") at `read_csv`

```
1 # data summary - overall
2 df |>
3   summary()
4
5 # structure of the data
6 glimpse(df)
```

- Why is `glimpse` especially useful having used the `col_types` with `read_csv`?

psych::headTail FUNCTION

- The `headTail` function from the *psyche* package

```
1 # data summary - head  
2 headTail(df)
```

summarytools PACKAGE

- We've seen how we can get an overview of the data using `dfSummary`
- This function is particularly useful for spotting anomalies in the data

```
1 dfSummary(df)
```

TIDYING VECTOR LABELS (I.E., COLUMN NAMES)

- We can use the `snakecase` package function `to_snake_case` to apply uniform *snake case* formatting to column names
- This addresses the point made earlier about consistent variable names

```
1 # check application of snake_case to data
2 to_snake_case(names(df))
3
4 # irreversibly apply snake_case to vector labels
5 names(df) <-
6   to_snake_case(names(df)) # lower_case with underlined separators
7
8 # examine the data
9 names(df)
```

CHANGING VECTOR NAMES

- You can also change the names of vectors using the `rename` function
- Note, the `rename` argument asks for the **new vector label followed by the old one**
 - this takes a little getting used to but its simple after enough practice

```
1 # Rename vectors to more descriptive labels
2 df |>
3   rename(perceived_robot_gender = alyx_is,
4         robot_expression = expression,
5         ppt_acc = resp_acc) |>
6   print()
```

- Using the code `names(df)` in the console, examine the vector labels. What does the output tell you?

RECODING VARIABLE LEVELS

- You may have noticed from the last session that there are some erroneous values in the sex vector.
- You can check this directly by using the `$` symbol to index the vector of interest in combination with the function `unique`.

```
1 # Check all unique entries in the vector `sex`  
2 unique(df$sex)  
3  
4 # Alternatively you can use tidyverse's `distinct`  
5 df |>  
6   select(sex) |>  
7   distinct()
```

- What does `select` do here?
- As you can see, something has gone wrong in the labelling of sex: there should only be values, “Female” and “Male”

USING `case_match` TO RECODE DATA

- Previously, the `recode` function was used in combination with `mutate` to recode values in a cell
- But more recently, people have moved to using `case_match`, as it is more intuitive and user friendly
- Info on `case_match`

case_match SYNTAX

```
1 # Using case_match to recode cell values
2 df <-
3   df |> # create new and tidy object for illustration
4   mutate(sex = case_match(sex,                      # mutate to change
5           "male" ~ "Male",                      # supply erroneous value and its
6           "ffemale" ~ "Female",                  # correct value
7           "femalew" ~ "Female",                  # preserve all other entries
8           .default = sex))                     # preserve all other entries
9
10 # Check the values of a vector using unique
11 unique(df$sex) # again, we use the $ symbol to isolate a specific vector
```

FOR LARGER DATASETS

- If you wanted to match values by a pattern, you can use `stringr` commands to recode variable levels
- In this case, we use the function `case_when` and `str_detect` to search the `sex` vector for values starting with either capital or lower case “m” and replace it with “Male”, and the same for principle “f”, to be recoded as “Female”

```
1 # Using case_when and str_detect to modify rows by string patterns
2 df |>
3   mutate(sex =
4     case_when(str_detect(sex,
5                 "^[Mm] ") ~ "Male",
6                 str_detect(sex,
7                 "^[Ff] ") ~ "Female",
8                 TRUE ~ sex))
```

DEALING WITH THE DATE COLUMN

- Dates are notoriously difficult to deal with in R and our date vector could not be read properly
- Thankfully there is a really useful package that simplifies everything; the aptly named `lubridate`
- It is already loaded as part of the Tidyverse
- Take a look at `lubridate` here

lubridate TO SORT OUT DATE COLUMN

- For more information see the [Tidyverse lubridate page](#)

```
1 df <-  
2   df |>  
3   mutate(exp_date = dmy(exp_date)) # mutate is a function that allows you  
4  
5 # Examine data  
6 glimpse(df)
```

- The date column should now be properly machine ready formatted (i.e., YYYY-MM-DD), based on the input formatting (i.e., dd/mm/yyyy)
- The machine readable format YYYY-MM-DD is preferable for almost all programming languages.
- In this format, you can ask R to make nuanced calculations between dates.

RECAP ON THE ASSIGN <- OPERATOR CONSIDERATIONS

- Irreversibly overwriting dataframe/tibble

```
1 df <-
2   df |>
3   select(id, eth)
4
5 # df only contains data about ppt id and eth
```

Create new dataframes/tibbles

```
1 df1 <-
2   df |>
3   select(id, eth)
4
5 # df1 only contains data about ppt id and eth
6 # df contains all 12 vectors
```

INTRODUCTION DATA WRANGLING

- You can create *tidyverse* scripts that summarise data similarly to **psych**
- **psych** is better for numeric data, which we'll be using next week.

```
1 df |>  
2   group_by(sex,  
3             resp_acc) |>  
4   summarise(n = n()) |>  
5   mutate(freq = n/sum(n)) |>  
6   mutate(freq = round(freq*100, digits = 2))
```

- Having had a look at the output, which is the above code doing? Think about each step of the process, punctuated by |>

EXERCISE

- What code would you write to calculate the proportion of different ethnicities and their prior experience with robots?
- How would we create a tibble called `df_non_native` that only included non-native English speakers in the data?

SAVE YOUR TIDY DATA

- Once you are satisfied with the amendments made to your data you can save it as a new tidy **.csv** file
- Again, use comments and bookmarks to signpost this part of your script

```
1 # save your new data object into `data_tidy`  
2 write_csv(df, # name of the object  
3           "data_tidy/robot-expression-data_tidy.csv") # desired directory
```

- Note: you can provide any name you like for the file in here, the important part is to pass **read_csv** the tibble of interest, and to include **data_tidy** in the stated directory

CLEANUP

```
1 # Clear data
2 rm(list = ls())  # Removes all objects from environment
3
4 # Clear packages
5 p_unload(all)  # Remove all contributed packages
6
7 # Clear plots
8 graphics.off()  # Clears plots, closes all graphics devices
9
10 # Clear console
11 cat("\014")  # Mimics ctrl+L
```

REFERENCES

