

Missing data, factors and dates in R

Nathan Green

Missing data

- As R was designed to analyse datasets, it includes the concept of missing data (which is uncommon in other programming languages)
- Missing data are represented in vectors as NA.
- When doing operations on numbers, most functions will return NA if the data you are working with include missing values.
- This feature makes it harder to overlook the cases where you are dealing with missing data.
- You can add the argument ***na.rm = TRUE*** to calculate the result while ignoring the missing values.

```
heights <- c(2, 4, 4, NA, 6)
mean(heights)
max(heights)
mean(heights, na.rm = TRUE)
max(heights, na.rm = TRUE)
```

Useful NA functions

If your data include missing values, you may want to become familiar with the functions `is.na()`, `na.omit()`, and `complete.cases()`. See below for examples.

```
## Extract those elements which are not missing values.  
heights[!is.na(heights)]  
  
## Returns the object with incomplete cases removed. The return  
na.omit(heights)  
  
## Extract those elements which are complete cases. The return  
heights[complete.cases(heights)]
```

Factors

- Factors are very useful and actually contribute to making R particularly well suited to working with data. So we are going to spend a little time introducing them.
- Factors represent categorical data. They are stored as integers associated with labels and they can be ordered or unordered.
- While factors look (and often behave) like character vectors, they are actually treated as integer vectors by R. So you need to be very careful when treating them as strings
- Once created, factors can only contain a pre-defined set of values, known as *levels*. By default, R always sorts levels in alphabetical order.

Examples

```
sex <- factor(c("male", "female", "female", "male"))
```

- R will assign 1 to the level "female" and 2 to the level "male" (because f comes before m, even though the first element in this vector is "male"). You can see this by using the function `levels()` and you can find the number of levels using `nlevels()` :

```
levels(sex)  
nlevels(sex)
```

Level order

- Sometimes, the order of the factors does not matter, other times you might want to specify the order because it is meaningful (e.g., “low”, “medium”, “high”), it improves your visualization, or it is required by a particular type of analysis. Here, one way to reorder our levels in the sex vector would be:

```
sex # current order
```

```
#> [1] male   female female male  
#> Levels: female male
```

```
sex <- factor(sex, levels = c("male", "female"))  
sex # after re-ordering
```

```
#> [1] male   female female male  
#> Levels: male female
```

Converting factors

- To convert to a character vector use

```
as.character(sex)
```

- In some cases, you may have to convert factors where the levels appear as numbers (such as concentration levels or years) to a numeric vector.
- eg, in one part of your analysis the years might need to be encoded as factors (e.g., comparing average weights across years) but in another part of your analysis they may need to be stored as numeric values (e.g., doing math operations on the years).
- `as.numeric()` returns the index values of the factor, not its levels, so it will result in an entirely new (and unwanted in this case) set of numbers.
- One method to avoid this is to convert factors to characters, and then to numbers.

Renaming factors

- When your data is stored as a factor, you can use the `plot()` function to get a quick glance at the number of observations represented by each factor level.
- Example

```
sex <- surveys$sex  
head(sex)
```

```
#> [1] M M  
#> Levels: F M
```

```
levels(sex)
```

```
#> [1] "" "F" "M"
```

```
levels(sex)[1] <- "undetermined"  
levels(sex)
```

```
#> [1] "undetermined" "F" "M"
```

```
head(sex)
```

```
#> [1] M M undetermined undetermined un  
determined  
#> [6] undetermined  
#> Levels: undetermined F M
```

Dates

- One of the most common issues that new (and experienced!) R users have is converting date and time information into a variable that is appropriate and usable during analyses.
- Simple practice for dealing with date data is to ensure that each component of your date is stored as a separate variable.
- Using `str()`, we can confirm that our data frame has a separate column for day, month, and year, and that each contains integer values

Lubridate

- We are going to use the `ymd()` function from the package `lubridate` (which belongs to the `tidyverse`; learn more [here](#)).
- `lubridate` gets installed as part as the `tidyverse` installation.
- When you load the `tidyverse` (`library(tidyverse)`), the core packages (the packages used in most data analyses) get loaded.
- `lubridate` however does not belong to the core `tidyverse`, so you have to load it explicitly with `library(lubridate)`

Date formats

- `ymd ()` takes a vector representing year, month, and day, and converts it to a Date vector.
- Date is a class of data recognized by R as being a date and can be manipulated as such.
- The argument that the function requires is flexible, but, as a best practice, is a character vector formatted as “YYYY-MM-DD”.
- Alternatives are YYYY/MM/DD or YYYY/DD/MM etc

Let's create a date object and inspect the structure:

```
my_date <- ymd("2015-01-01")  
str(my_date)
```

Now let's paste the year, month, and day separately - we get the same result:

```
# sep indicates the character to use to separate each component  
my_date <- ymd(paste("2015", "1", "1", sep = "-"))  
str(my_date)
```