

Problem Set 9

Andreas Bender, Philipp Kopper, Sven Lorenz

25 January 2022

Read chapter 13 on relational data in R4DS. You can skip the exercises.

Exercise 1

In this exercise we investigate data sets on baseball. In particular, we are interested in creating player statistics. You will combine different data sets to gather all the information that you are looking for. Save your resulting objects with the same naming convention as last week. That means your solution for exercise 1a) should be saved as `ex1a.Rds` in the folder `solutions`. Use the `all.equal` function to check your solutions with the solutions available from Moodle.

- a) Load the `Lahman` package and have a look at the `Pitching` dataset. Use the `?` function to get an overview on the following data sets: `Pitching`, `People`. Filter the `Pitching` data so that it only contains the data from the years 2018 and 2019. What is the primary key for this data set? Work with this filtered data set throughout the problem set.
- b) Group the data by player, year and team. Count the number of wins, losses, games, and home runs for each player, year and team. Store the resulting data in an object called `pitching`.
- c) What is the primary key of the `pitching` data set?
- d) Have a look at the `People` data set from the `Lahman` package. What is the primary key here?
- e) For each player in your `pitching` data set, we now want to have information on their birthplace. Join the appropriate data to your pitching table. Keep all rows from the pitching data and all columns from both tables.
- f) We now want to summarize the pitching information for each player regardless of year and team. Group the `pitching` data by `playerID` and add up the information about wins, losses, games and home runs.
- g) Now join this summarised pitching data set with the people data set as in task e). Based on this data set, create a data frame that for each country contains the information about the number of players from the respective country (excluding the USA). Sort the data such that the country with most players is at the top and the country with the least players is at the bottom. Use `ggplot` to create a bar plot that shows the number of players per country (Hint: use `geom_col`). Make sure that the country names in the axis labels don't overlap. **Bonus:** Recreate the above plot such that the bars are ordered from highest to lowest.

Exercise 2

We collected data on the coffee drinking preferences in the office. The people working from the office buy a variation of beans of which everyone can take some. We stored these information in the object `people`. The `data.frame` consists of the Name of the consumer, the information how many cups they drink, where the coffee is from (which is in our scenario identical to its name) and what kind of coffee drink they prefer. Note

that there are two special cases: Julia drinks her own coffee of which we do not know a lot and Martin does not drink coffee at all. `coffee` features information on the beans that the people drink, precisely on the roast (light to dark) and its price (per kilogram). Next to these, we have a `data.frame` that explains how many shots of espresso we need for each of the drinks (`drinks`). Additionally, it tells how much milk (in ml) one needs for the drinks. Note that for one espresso shot you typically need 18g of coffee. A litre of milk costs 1 euro.

```
people <- data.frame(
  name = c("Andreas", "Philipp", "Sven", "Martin",
           "Moritz", "Julia", "Jana", "Klaus",
           "Maria", "Leonie", "Frederick", "Ute"),
  n_cups = c(3, 5, 1, NA,
            3, 4, 1, 1,
            6, 5, 2, 1),
  origin = c("Ethiopia", "Honduras", "Costa Rica", NA,
            "Vietnam", "Brazil", "Honduras", "Honduras",
            "Costa Rica", "Vietnam", "Colombia", "Costa Rica"),
  drink = c("Cappuccino", "Flat White", "Cappuccino", NA,
            "Espresso", "Latte", "Cappuccino", "Cappuccino",
            "Espresso doppio", "Flat White", "Cappuccino", "Espresso doppio"))

coffee <- data.frame(
  origin = c("Ethiopia", "Honduras", "Costa Rica", "Vietnam", "Colombia"),
  roast = c("decaf", "medium", "dark", "light", "medium"),
  price = c(18.95, 26.15, 17.99, 29.00, 21.40))

drinks <- data.frame(
  drink = c("Espresso", "Espresso doppio", "Cappuccino", "Cafè Crema",
            "Latte", "Flat White", "Americano"),
  shots = c(1, 2, 1, 1, 1, 2, 1),
  milk = c(0, 0, 100, 0, 150, 100, 0)
)
```

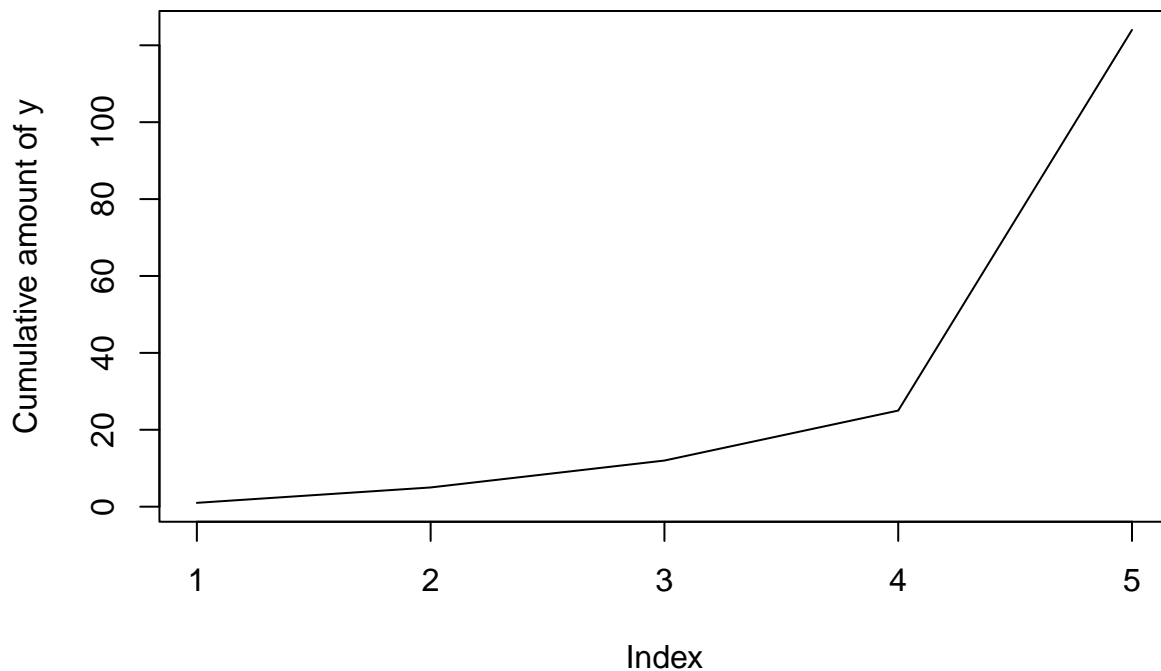
- Compute how many shots of each coffee are pulled on a regular office day. How much do the colleagues spend per day on coffee beans? How much for milk?
- Compute how many kilograms of each roast of coffee are consumed over the course of the week.
- Identify the individual costs of coffee consumption per day. Note that Julia pays for her own coffee, so she has an individual contribution of 0 Euros. As most people have the feeling that the distribution of coffee consumption is extremely uneven, they want to express this inequality by a number, namely the Gini coefficient. Use the `ineq` package to compute the Gini coefficient (`ineq` function) of this distribution. Interpret the value. Also draw the Lorenz curve. You find a comprehensive overview on the Gini coefficient and the Lorenz curve here, if you did not attend Statistics 1 this semester or as a refresher. Note that the package does most of the work for you in computing the Gini coefficient and drawing the curve. The Gini coefficient can be directly inferred from the Lorenz Curve. To draw the Lorenz curve one needs to sort all values of a vector y from lowest to highest. Then one computes the cumulative value of y over the indices. To plot the curve one displays the index on the x-axis and the cumulative value of y on the y-axis. In the case of a uniform distribution, the curve equals a diagonal line. We illustrate the procedure of computing the respective y-axis values and plotting below (Note that often times, the Lorenz curve is plotted not using the absolute numbers but the relative shares on both axes.):

```
y <- c(99, 4, 7, 1, 13)
df <- data.frame(index = 1:length(y), y = sort(y), y_cumu = cumsum(sort(y)))
df
```

```
##   index  y y_cumu
## 1     1  1     1
## 2     2  4     5
## 3     3  7    12
## 4     4 13    25
## 5     5 99   124
```

```
plot(df$index, df$y_cumu, main = "Example for Lorenz Curve", type = "l",
     xlab = "Index",
     ylab = "Cumulative amount of y")
```

Example for Lorenz Curve



this will both be cared for by the package that you are supposed to use.

Session Info

```
sessionInfo()
```

```
## R version 4.1.0 (2021-05-18)
## Platform: x86_64-apple-darwin17.0 (64-bit)
## Running under: macOS Big Sur 10.16
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.1/Resources/lib/libRblas.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.1/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
```

```
## [1] stats      graphics  grDevices utils      datasets  methods   base
##
## loaded via a namespace (and not attached):
## [1] compiler_4.1.0    magrittr_2.0.1    tools_4.1.0       htmltools_0.5.1.1
## [5] yaml_2.2.1        stringi_1.7.3     rmarkdown_2.11    highr_0.9
## [9] knitr_1.33        stringr_1.4.0     xfun_0.24         digest_0.6.27
## [13] rlang_0.4.11      evaluate_0.14
```