

Data Handling: Import, Cleaning and Visualisation

Exercise to lecture 4: csv and arrays

Dr. Aurélien Sallin

Working with a data frame

Set your script

Set your R Script.

```
#####  
# Data Handling Course: Example Script for Data Gathering and Import  
#  
# Imports data from ...  
# Input: import c to data sources (data comes in ... format)  
# Output: cleaned data as CSV  
#  
# A. Sallin, St. Gallen, 2023  
#####  
  
# SET UP -----  
# Load packages  
library(readr)  
  
# SET PATH -----  
# If not in NUVOLOS, set correct path!  
# financial_data <- read.csv("Path/to/my/file/financial_data.txt")
```

Import data

Have a look at the file `financial_data.txt` using your favorite text editor. What do you notice?

Import the table using the `read.csv()` function in your environment. Make sure you have the right path to access the `.txt` document. What does this parser do? Explore the `data.frame`. What is its structure? What are its dimensions?


```
# The data looks now like what I expect... except for the variable "Revenue", which
# is a character. I remove the special ö character from the encoding issue and
# coerce
head(financial_data, 10)
```

```
##      Firm Year Revenue    Profit Category
## 1 FirmA 2017    4355    897.4552 Finance
## 2 FirmB 2017    4919   1091.3730 Health
## 3 FirmC 2017    4065   1231.3810 Tech
## 4 FirmD 2017    4989    860.2956 Tech
## 5 FirmE 2017    4172   1684.9384 Tech
## 6 FirmA 2018    2003    361.5399 Tech
## 7 FirmB 2018    1622    330.1158 Health
## 8 FirmC 2018    3952   1963.5914 Tech
## 9 FirmD 2018    3692   1561.4979 Finance
## 10 FirmE 2018   1933ö    621.1375 Finance
```

```
str(financial_data)
```

```
## 'data.frame':    30 obs. of  5 variables:
## $ Firm      : chr  "FirmA" "FirmB" "FirmC" "FirmD" ...
## $ Year      : int   2017  2017  2017  2017  2017  2018  2018  2018  2018  2018 ...
## $ Revenue   : chr   "4355" "4919" "4065" "4989" ...
## $ Profit    : num   897  1091  1231  860  1685 ...
## $ Category  : chr   "Finance" "Health" "Tech" "Tech" ...
```

```
financial_data[10, 3] <- 1933

# Coerce to numeric
financial_data$Revenue <- as.numeric(financial_data$Revenue)

# Another way of writing the column selection
financial_data[10, "Revenue"]
```

```
## [1] 1933
```

```
financial_data[10, "Revenue"] <- 1933

financial_data[, "Revenue"] <- as.numeric(financial_data[, "Revenue"])

# The data is now ready. You are ready to compute the rest of the exercise.
# END - for now
```

Summary statistics of your data

Compute the summary statistics for each variable using the `summary()` command. What does this command give you? What do you notice? Make the necessary changes.

```
##      Firm              Year      Revenue      Profit
## Length:30      Min.    :2017      Min.    :1269      Min.    : 164.4
## Class :character 1st Qu.:2018      1st Qu.:2332      1st Qu.: 665.5
## Mode  :character Median :2020      Median :3610      Median : 990.4
##              Mean  :2020      Mean  :3317      Mean  :1001.0
##              3rd Qu.:2021      3rd Qu.:4048      3rd Qu.:1245.4
##              Max.   :2022      Max.   :4989      Max.   :1963.6
##      Category
## Length:30
## Class :character
## Mode  :character
##
##
##
```

Variable creation

Create a new variable “costs”, which is the revenue - profit. [There are many ways to create a variable in a data frame. Here, use the `$` index.]

Factor variable

Which variable is (should be) a factor? Recode this variable as a factor. What are the levels? Should we have the variable `Firm` as a factor?

Nests - more difficult question... but still exam relevant



Split your data using the factor variable into three data frames that are contained in a list. Compute the mean profit for each data frame.

- Hint: use the function `split`.
- Hint: use a `for`-loop over each list element to compute the mean

Advanced: map (not exam relevant)

Do the same as the exercise above using the `map` function. Install the packages `tidyr`, `dplyr`, and `purrr`.