# **Data Analysis with R:**
## Lecture Slides: Day 1 - Monday

Sonja Hartnack, Terence Odoch & Muriel Buri

22nd of July 2019

**Goals of the course**

To be able to...

- import data sets to R

- describe data with R

- apply basic statistical tests in R

- some ideas for more advanced statistical tools ...

- simulate a data set similar to own research

**General remarks**

**Course schedule:**

- Starting at 9:00am / 9:30am (?)
- Tea breaks in between
- Lunch break
- Teaching until 4.30pm ($\sim$ 5pm)

**Obtaining a certificate is conditional on:**

- active participation in class
- attending at least 75 % of the course (lecture & exercises)
- assignments during now and October
- short final exam in October (format to be defined)

**Getting to know each other**

- My name is ...

- I am doing a Master / a PhD in ...

- I hope to learn in this course how to ....

- My personal goal for this course is ...

**How do we reach these goals**

- hands on exercises with R:
    - `chickwts`
    - `ToothGrowth`
    - `bacteria`
    - `perulung`
    - ... and others.

- interactive discussions & student's present their own solutions

- ask us a lot of questions but also ask google for help!

- group work

- short motivational lectures

# Do you all have RStudio and R installed on your computers?

# Get started with data set: chickwts

An experiment was conducted to measure and compare the effectiveness of various feed supplements on the growth rate of chickens.

```r
# load data set "chickwts"
data("chickwts", package = "datasets")
# the head(...) function shows the first 6 observations
head(chickwts)

##   weight      feed
## 1    179 horsebean
## 2    160 horsebean
## 3    136 horsebean
## 4    227 horsebean
## 5    217 horsebean
## 6    168 horsebean

# FUNCTION - open bracket - DATA SET / VARIABLE - close bracket
```
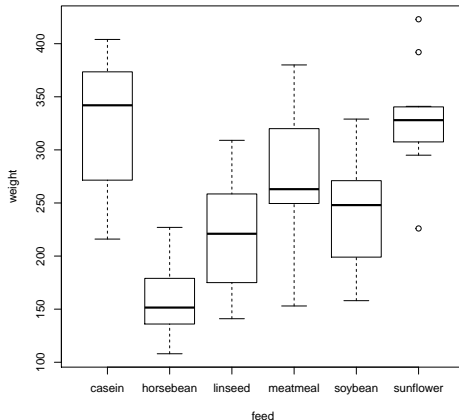
# Ideas for plotting the data

# Ideas for plotting the data

```
# use x axis to show the categorical variable (feed),
# y axis to represent the continuous variable (weight)
# boxplot (y.cont.variable ~ x.cat.variable, data = dataset)
# ?boxplot
boxplot(weight ~ feed, data = chickwts)
```
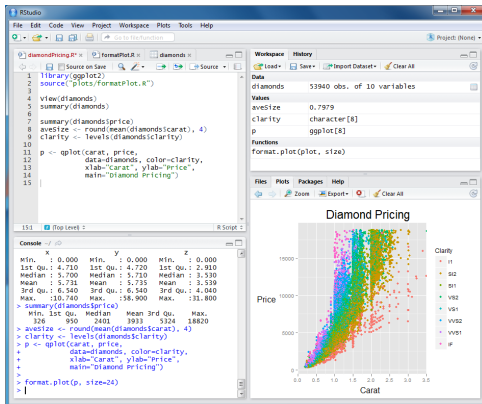
# Exercise: Statistical Terminologies

## Functionalities in RStudio

- Source
- Console
- Environment, History, Files
- Files, Plots, Packages, Help

# Good housekeeping!

- Define manually a new folder called **rcourse** in your personal documents on your personal computer

- Know in which directory you are

```
getwd()

## [1] "/home/mburi/Documents/git_svn/DataAnalysisWithR/Lectures"
```

- Set directory path

```
# back- and forslash is dependent on the system
setwd("C:/Users/muriel/Documents/rcourse/")
setwd("C:\\Users\\muriel\\Documents\\rcourse\\")
```

- Always clean up before starting with new R-Script

```
rm(list=ls()) # empty workspace, delete previously saved variables
```

# How to get help in R

```
?chickwts
?boxplot
```

Also, have a look at the examples at the end of the help pages.

# Exercise: Getting to know R and `chickwts`

# A data frame in R: `chickwts`



**chickwts[ ROWS , COLUMNS ]**

chickwts[ 6, 1 ]

chickwts[ 11, 2 ]

## Rows and columns of a data frame: `chickwts`

Values of ...

```r
# Load (internal) data set from R
data("chickwts")

# ... all columns of sixth observation:
chickwts[6, ]

# ... all columns of sixth to eleventh observation:
chickwts[c(6:11), ]

# ... all columns of sixth, eleventh and twentieth observation:
chickwts[c(6, 11, 20), ]

# ... all rows of first column (weight):
chickwts[ , 1]

# ... all rows of second column (feed):
chickwts[ , 2]

# or use the "$" sign as a reference to column "feed":
chickwts$feed
```

# Exercise: Summary Statistics for the `chickwts` data set

# Rules for importing data into R

- First row of excel sheet contains **variable names**:
  `y`, `ap`, `hilo`, `week`, `ID`, `trt`.
- Columns of excel sheet represent **variables**.
- Rows of excel sheet represent **observations per individual** (except for the first row).
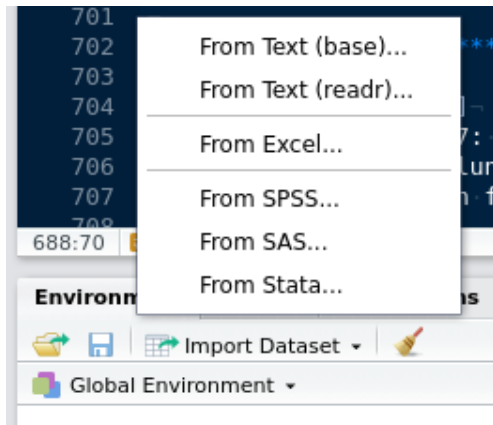
## Rules for naming variables

Variable names should ...

- start with a letter (not a number):
  `y, ap, hilo, week, ID, trt`

- longer variables names should be separated with dots:
  `time.in.weeks`

- do not use special characters, such as /, #, @, &, $\star$, ...

**How to import external data files into R?**

- > Import Dataset > **From Text (base)...** > CSV Files (.csv) or

- Environment (upper right corner)

- > Import Dataset > **From Text (base)...** > CSV Files (.csv)

```
perulung_ems <- read.csv("perulung_ems.csv", row.names = 1,
                          sep = ";")
lung <- data.frame(perulung_ems)
```

- > Import Dataset > **From Text (base)...** > Text Files (.txt)
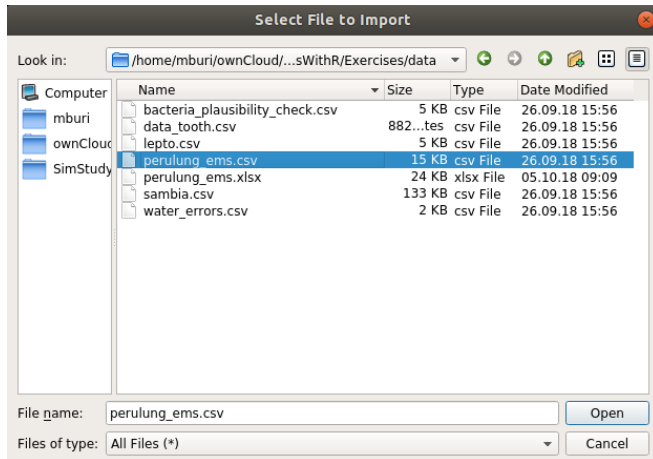
## How to import .txt and .csv files into R? (2/3)

- Environment (upper right corner)
- $>$ Import Dataset $>$ From Text (base)... $>$ CSV Files (.csv)

# How to import .txt and .csv files into R? (2/3)

- Environment (upper right corner)
- \> Import Dataset \> From Text (base)... \> CSV Files (.csv)

## How to import .txt and .csv files into R? (3/3)



```
perulung_ems <- read.csv("perulung_ems.csv", row.names = 1,
                         sep = ";")
lung <- data.frame(perulung_ems)
```

**Exercise: Data import to R and summary statistics** `perulung_ems.csv` **R**

Data from a study of lung function among children living in a deprived suburb of Lima, Peru. Data taken from Kirkwood and Sterne, 2nd edition.

Variables:

- `fev1`: in liter, "forced expiratory volume in 1 second" measured by a spirometer. This is the maximum volume of air which the children could breath out in 1 second
- `age`: in years
- `height`: in cm
- `sex`: 0 = girl, 1 = boy
- `respsymp`: respiratory symptoms experienced by the child over the previous 12 months