# Output development using R & R markdown

Ziv Shkedy and Thi Huyen Nguyen

Hasselt University

Foundations for inference

Ha Noi, 29/05/2024

ER-BioStat

https://github.com/eR-Biostat

@erbiostat

# The `airquality` data

## Introduction

# What do we do in this session ?

- We conduct a simple analysis for the variable Wind speed in the `airquality` data:

  - Summary statistics.

  - Graphical display: histogram.

  - Confidence interval.

  - Test of hypothesis.

- Focus:

  - How to produce an output ?

  - How to combine text and software output in the same document ?

# The `airquality` data

## Part 1

### The dataset

# The `airquality` data in R

```
> dim(airquality)
[1] 153 6
> names(airquality)
 [1] "Ozone" "Solar.R" "Wind" "Temp" "Month" "Day"
```

The R object for the data: 153 observations and 6 variables.
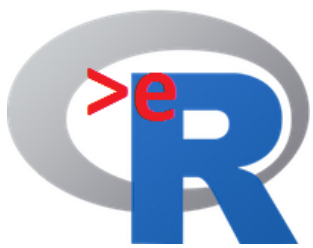
Variables names:

Ozone: Mean ozone in parts per billion from 1300 to 1500 hours at Roosevelt Island.
Wind: Average wind speed in miles per hour at 0700 and 1000 hours at LaGuardia Airport.

# The `airquality` data in R

```
> head(airquality)

    zone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
5    NA      NA 14.3   56     5   5
6    28      NA 14.9   66     5   6
```

# The `airquality` data

## Part 2

## The Rmd program

# The output

- We run the R markdown file.
- Produce output format: Word document.

# The Rmd program



```
1   ---
2   title: 'Case study 1: analysis of the daily average wind speed in New York in 1973'
3   output:
4     word_document: default
5     pdf_document: default
6     html_document: default
7   subtitle: Foundations for inference using R
8   layout: page
9   ---
10
11
12  ```{r setup, include=FALSE}
13  options(htmltools.dir.version = FALSE)
14  knitr::opts_chunk$set(echo = TRUE,
15                        message = FALSE,
16                        warning = FALSE,
17                        eval = TRUE,
18                        tidy = FALSE)
19  library(knitr)
20  library(tidyverse)
21  library(deSolve)
22  library(minpack.lm)
23  library(ggpubr)
24  library(readxl)
25  library(gamlss)
26  library(data.table)
27  library(grid)
28  library(png)
29  library(nlme)
30  library(gridExtra)
31  library(mvtnorm)
32  library(e1071)
33  library(lattice)
34  library(ggplot2)
```

Document setup.

Many R packages, not all needed.

# The Rmd program

# Choose the output

# The Word document output

# The `airquality` data

## Part 3

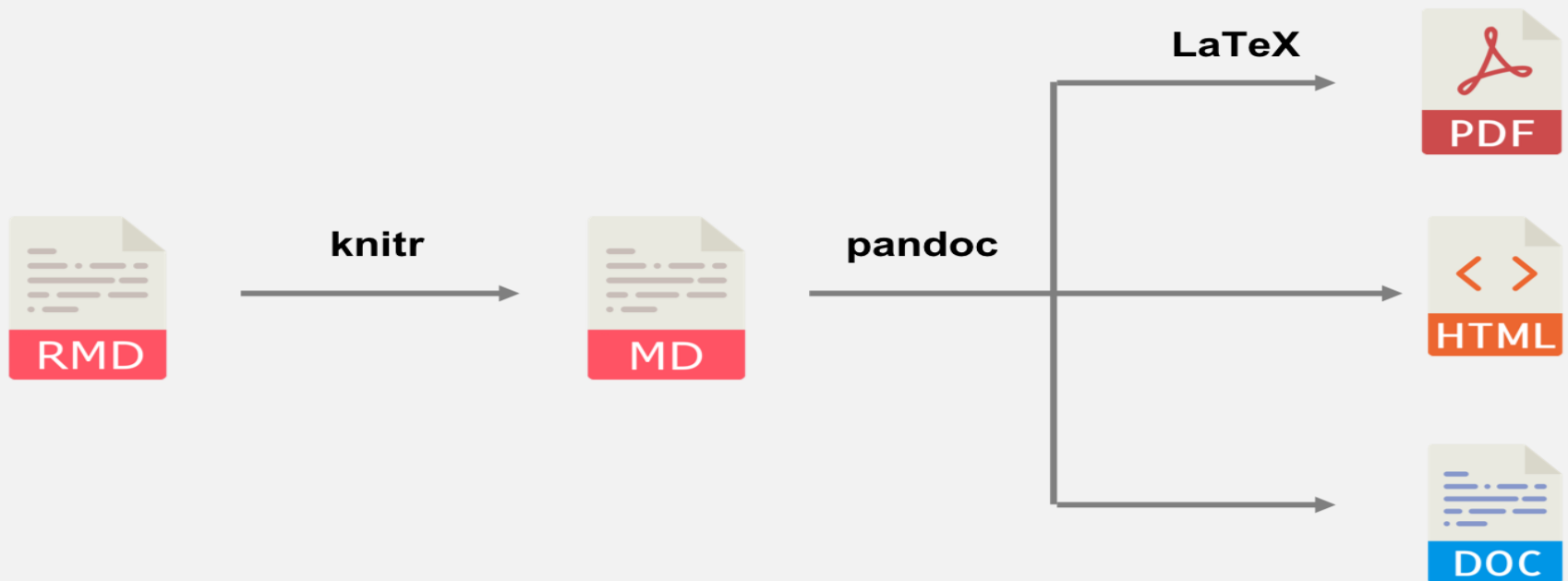### How to produce the Word output ?

# Reproducible Research

- Aim: create an output in a Word document.

- Can be used to communicate the analysis' results with other people in the organization.

- Not all potential readers are interested on "how to do the analysis".

- We DO NOT aim to develop a report for the analysis but to provide a document from which the results can be seen and discuss by different people in the organization.

# The Rmd file

- Analyses ⟶ high quality report.
- Rmarkdown – Different dynamic and statistic formats (html, pdf, word, books, dashboard, e.t.c).

# The Word output



- A Word document output.
- Presents the same analysis as in the example.

# Part 3.1:
# How to set up the Word file?

# The Rmd file

- We use Rmd file to
  - Conduct the analysis.
  - Set up the document.

- We use a Word file in order to
  - Present & communicate the result.

# Set up the document

# Titles, authors and dates

# Titles, authors and dates



**Case study 1: analysis of the daily average wind speed in New York in 1973** — Title

**Foundations for inference using R** — Sub title
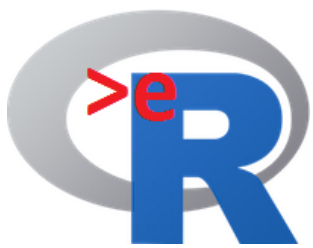
Ziv Shkedy et al.

29-05-2024

## Case study 1: The wind speed in the airquality dataset

### Exploratory analysis of the daily average of the wind speed

The airquality dataset is a R object gives information about 153 daily air quality measurements ($n = 153$) in New York, May to September 1973.

```
dim(airquality)

## [1] 153    6

names(airquality)

## [1] "Ozone"   "Solar.R" "Wind"    "Temp"    "Month"   "Day"

head(airquality)

##    Ozone Solar.R Wind Temp Month Day
```

# Part 3.2:
# The Word document and the Rmd program in details.

# Section, subsection, subsubsection



- In addition to the code, we can add free text in the Rmd file.

# Section, subsection, subsubsection



The slide shows a Word document screenshot with the following visible content:

**Case study 1: analysis of the daily average wind speed in New York in 1973**

**Foundations for inference using R**

Ziv Shkedy et al.

29-05-2024

**Case study 1: The wind speed in the airquality dataset** ← section

**Exploratory analysis of the daily average of the wind speed** ← subsection

free text →

The airquality dataset is a R object gives information about 153 daily air quality measurements ($n = 153$) in New York, May to September 1973.

```
dim(airquality)
## [1] 153   6
names(airquality)
## [1] "Ozone"   "Solar.R" "Wind"   "Temp"   "Month"   "Day"
head(airquality)
##   Ozone Solar.R Wind Temp Month Day
```

The code is shown as a part of the output →

24

# Code in the Rmd file



The code for histogram

# The output in the Word file



Histogram with density of wind speed.

# Code in the Rmd file

# The output in the Word file

# Code in the Rmd file

# The output in the Word file



**Test of hypothesis about the population mean**

Testing the hypotheses whether the wind speed is equal to 9 versus a two-sided alternative hypothesis at the significant level of 0.05 can be formulated by:

$$H_0: \mu = 9 \quad \text{Vs.} \quad H_1: \mu \neq 9.$$

We use the z.test() function and specify mu=9$.

```
z.test(wind, SD.wind, mu=9)
```

```
##
##  One Sample z-test
##
## data:  wind
## z = 3.3619, n = 153.00000, Std. Dev. = 3.52300, Std. Dev. of the sample
## mean = 0.28482, p-value = 0.0007742
## alternative hypothesis: true mean is not equal to 9
## 95 percent confidence interval:
##    9.399284 10.515749
## sample estimates:
## mean of wind
##     9.957516
```

Since p-value = 0.0007742 which is much smaller than $\alpha = 0.05$, there is sufficient evidence to say that the mean of the wind speed is not equal to 9.

The output for testing hypotheses.

# Discussion

- R Studio + R markdown:

- Easy to use.

- Text + code.

- Output:
  - Standard: HTML, PDF, DOC.
  - Example: Word.doc.