# Reproducible Research

Şebnem Er

University of Cape Town, Statistical Sciences Department

Sebnem.Er@uct.ac.za

December 05, 2017

| | | | |
|---|---|---|---|
| _ER 14 May.doc | 2017/02/07 3:02 P... | Microsoft Word 9... | 502 KB |
| _ER 14 May_.doc | 2017/12/03 12:16 ... | Microsoft Word 9... | 502 KB |
| x_ER 6 Dec16.docx | 2017/03/26 3:20 P... | Microsoft Word D... | 162 KB |
| _ER.docx | 2012/10/27 12:57 ... | Microsoft Word D... | 56 KB |
| x_ERtoday.docx | 2013/10/10 9:16 A... | Microsoft Word D... | 77 KB |

$$\bar{Y}_h = \sum_{i=1}^{N_h} y_{hi} / N_h$$

Population mean of elements in stratum h

$$\bar{y}_h = \sum_{i=1}^{n_h} y_{hi} / n_h$$

Sample mean of elements in stratum h

$$S^2 \quad \sum_{i=1}^{N_h} (y_{hi} - \bar{Y}_h)^2 /$$

Population variance of the elements in stratum h

$$\bar{Y}_h = \sum_{i=1}^{N_h} y_{hi}/N_h$$

Population mean of elements in stratum h

$$\bar{y}_h = \sum_{i=1}^{n_h} y_{hi}/n_h$$

Sample mean of elements in stratum h

$$\sigma_h^2 = \sum_{i=1}^{N_h} \left(y_{hi} - \bar{Y}_h\right)^2/(N_h - 1)$$

Population variance of the elements in stratum h

Hi ⬚

How are you? Sorry for my late reply. Usual excuses on my side, teaching taking over every little time.

I worked a bit on the paper. I am attaching it here. I exclude the multivariate part for now. All the references etc. need adjustments. Will get to that. I started the whole work on Latex because Word is causing a bit too much of a trouble for formula. I attach both the pdf and the tex files. Feel free to comment on either.

Best wishes,
Sebnem

**2 Attachments**

📄 **Manus.pdf**

📄 **Manus.tex**

Hi Sebnem

Thanks. But we've got a problem. I don't work on LaTeX, at all. What to do know?

Best,

**Sebnem Er** <er.sebnem@gmail.com>                                    Mar 31

to

I am happy to do all the typing in Latex, it is very similar to Word. You can make your comments in the tex file. All you need to do is typing your comments or additions starting with a % (percentage sign) and save it and send me the file, I will compile it.
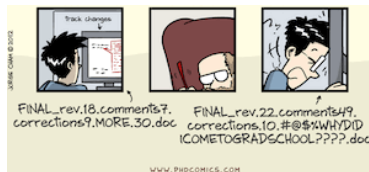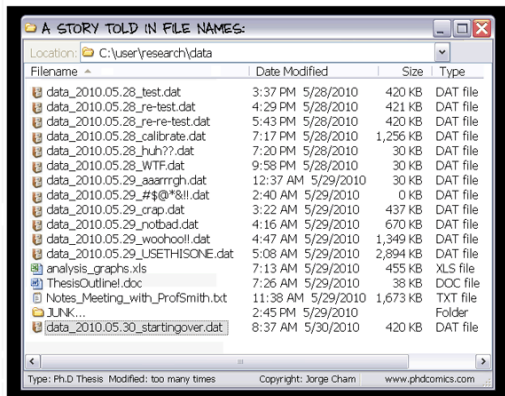
It is a real mission for me to change all the formula (which convert themselves to images) when the file is transfered between you and me. Incompatible word problems. The formulas get repeated.

Just have a look at the tex file (you can open it with . It is similar to an R file so actually you can open it in R and edit it there.

Let me know if it is OK.

```r
```{r getting_the_dataset, echo=FALSE}
#Martin_SPSS_raw_data_original <- read_sav("C:/Users/01438475/Google Drive/Research/Chao/M Schade/Martin
SPSS raw data originalmissing2.sav")

Martin_SPSS_raw_data_original <- read.csv("C:/Users/01438475/Google Drive/Research/Chao/M Schade/Martin
Excel raw data - original 2017nomissing.csv")

martin = as.data.frame(Martin_SPSS_raw_data_original)
martin$Gender = factor(martin$Gender)
levels(martin$Gender)[levels(martin$Gender)==1] <- "Male"
levels(martin$Gender)[levels(martin$Gender)==2]  <- "Female"
```
```

# A genetic algorithm approach to determine stratum boundaries and sample sizes of each stratum in stratified sampling

Timur Keskintürk ✉, Şebnem Er ✉

⊞ **Show more**

Get rights and content

For the comparison we examined eight examples (data can be obtained from URL http://www.isletme.istanbul.edu.tr/ogrelem/timur/English/stratification.htm) with different characteristics. The first example (iso487) consists of 487 Turkish

In all of these examples, you will see problems from

- Researchers/Collaborators
- Students
- Instructors
- Editors
- Private Sector Researchers
- Others who want to access your data and repeat what you have done

point of view. Everyone is affected.

Science, according to the American Physical Society, "is the systematic enterprise of gathering knowledge, organizing and condensing that knowledge into testable laws and theories."

Science, according to the American Physical Society, "is the systematic enterprise of gathering knowledge, organizing and condensing that knowledge into testable laws and theories."

How do we evaluate scientific claims?

Science, according to the American Physical Society, "is the systematic enterprise of gathering knowledge, organizing and condensing that knowledge into testable laws and theories."

How do we evaluate scientific claims?

Replication

Science, according to the American Physical Society, "is the systematic enterprise of gathering knowledge, organizing and condensing that knowledge into testable laws and theories."

How do we evaluate scientific claims?

Replication

Same results again and again $->$ findings are relevant

Science, according to the American Physical Society, "is the systematic enterprise of gathering knowledge, organizing and condensing that knowledge into testable laws and theories."

How do we evaluate scientific claims?

Replication

Same results again and again $->$ findings are relevant

However, replication "requires the complete and open exchange of data, procedures, and materials". We cannot replicate many scientific research, due to time, money or uniqueness of the research.

Then what do we do? We can have a middle ground and deploy reproducible research instead, which is the calculation of quantitative scientific results by independent scientists using the original datasets and methods.

Then what do we do? We can have a middle ground and deploy reproducible research instead, which is the calculation of quantitative scientific results by independent scientists using the original datasets and methods.

- What are the advantages of Reproducible Research?

Then what do we do? We can have a middle ground and deploy reproducible research instead, which is the calculation of quantitative scientific results by independent scientists using the original datasets and methods.

- What are the advantages of Reproducible Research?
  - Better work habits
  - Better team work
  - Changes are easier
  - Higher research impact

Then what do we do? We can have a middle ground and deploy reproducible research instead, which is the calculation of quantitative scientific results by independent scientists using the original datasets and methods.

- What are the advantages of Reproducible Research?
  - Better work habits
  - Better team work
  - Changes are easier
  - Higher research impact

- What tools can we use for RR?

- Before explaning the tools used for RR, what is the life cycle of a research?
  - Data collection
  - Data cleaning
  - Statistical analysis
  - Presentation of results/documentation/editing

For every step of your research, you need a tool for RR. Everything in one script.

1. Document everything!,
2. Everything in a script file,
3. All files should be human readable: Literate Programming,
4. Explicitly tie your files together: data, codes, pdf output etc.,
5. Have a plan to organize, store, and make your files available.

Using these tips will help make your computational research really reproducible.

# 1. Document everything

Ideally, you should tell your readers how you

- gathered your data,
- analyzed it, and
- presented the results.
- a key part of documenting with R is that you should record your session info:

```
sessionInfo()
```

# 2. Everything is a script file

- .R file
- .txt file
- .Rmd file
- .m file

etc.

# 3. All files should be human readable

Treat all of your research files as if someone who has not worked on the project will, in the future, try to understand them.

Treat all of your research files as if someone who has not worked on the project will, in the future, try to understand them.

**including yourself!**

# 3. All files should be human readable

Treat all of your research files as if someone who has not worked on the project will, in the future, try to understand them.

**including yourself!**

With this in mind it is a good idea to comment frequently.

# 3. All files should be human readable

Treat all of your research files as if someone who has not worked on the project will, in the future, try to understand them.

**including yourself!**

With this in mind it is a good idea to comment frequently.

*Commenting Guidelines*

- write a comment before a block of code describing what the code does,
- comment on any line of code that is ambiguous

- **R**: First step is to prepare your research in an environment where you can type your code and the machine can convert the code into analysis, in between you can leave comments what each code is doing.

# Tools for Reproducible Research (RR)?

- **R**: First step is to prepare your research in an environment where you can type your code and the machine can convert the code into analysis, in between you can leave comments what each code is doing.
- **RStudio**: RStudio allows you to do all of the things R can do. It is a happy medium between R's text-based interface and a pure GUI and it can be linked to numerous reproducible research publishing environments such as LaTeX.

- **R**: First step is to prepare your research in an environment where you can type your code and the machine can convert the code into analysis, in between you can leave comments what each code is doing.
- **RStudio**: RStudio allows you to do all of the things R can do. It is a happy medium between R's text-based interface and a pure GUI and it can be linked to numerous reproducible research publishing environments such as LaTeX.
- **knitr**: an R package for literate programming, i.e. it allows you to combine your statistical analysis and the presentation of the results into one document. Yihui Xie is the developer. see: https://yihui.name/knitr/

- **R Markdown**: One of the document formats that knitr supports, and it is also the simplest one. Markdown is a both easy-to-read and easy-to-write language.

- **R Markdown**: One of the document formats that knitr supports, and it is also the simplest one. Markdown is a both easy-to-read and easy-to-write language.
- **Cloud storage & versioning**: Services such as Dropbox and Git/Github that can store data, code, and presentation files, save previous versions of these files, and make this information widely available.

- **R Markdown**: One of the document formats that knitr supports, and it is also the simplest one. Markdown is a both easy-to-read and easy-to-write language.
- **Cloud storage & versioning**: Services such as Dropbox and Git/Github that can store data, code, and presentation files, save previous versions of these files, and make this information widely available.
- **Unix-like shell programs**: These tools are useful for working with large research projects.

# Knitr and R Markdown

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.
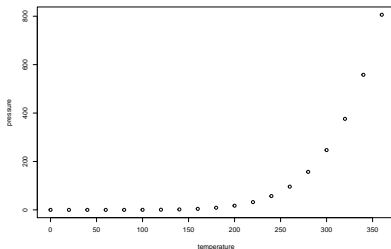
```
plot(pressure)
```



Figure 1: Scatterplot of Pressure vs Temperature

```
# Fit simple linear regression model
M1 <- lm(Examination ~ Education, data = swiss)
```

% latex table generated in R 3.4.2 by xtable 1.8-2 package % Tue
Dec 05 18:08:46 2017

Table 1: Linear Regression, Dependent Variable: Exam Score

|             | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 10.1     | 1.3        | 7.9     | 0.0      |
| Education   | 0.6      | 0.1        | 6.5     | 0.0      |

This was created with **xtable** package. For multiple tables, you can
use **apsrtable** package in **R**.

We need to be able to access our files from multiple devices in different locations.

# Cloud Storage and Versioning

We need to be able to access our files from multiple devices in different locations.

We often need a way for our collaborators to access and edit research files as well.

We need to be able to access our files from multiple devices in different locations.

We often need a way for our collaborators to access and edit research files as well.

When working on a collaborative project, one of the authors may accidentally delete something in a file that another author needed.

We need to be able to access our files from multiple devices in different locations.

We often need a way for our collaborators to access and edit research files as well.

When working on a collaborative project, one of the authors may accidentally delete something in a file that another author needed.

To deal with these issues we need to store our data in a system that has version control.

## Cloud Storage and Versioning

We need to be able to access our files from multiple devices in different locations.

We often need a way for our collaborators to access and edit research files as well.

When working on a collaborative project, one of the authors may accidentally delete something in a file that another author needed.

To deal with these issues we need to store our data in a system that has version control.

Version control systems keep track of changes we make to our files and allows us to access previous versions if we want to.

At the heart of GitHub is an open source version control system (VCS) called Git. Git is responsible for everything GitHub-related that happens locally on your computer.

At the heart of GitHub is an open source version control system (VCS) called Git. Git is responsible for everything GitHub-related that happens locally on your computer.

To use Git on the command line, you'll need to download, install, and configure Git on your computer.

# Git/Github

At the heart of GitHub is an open source version control system (VCS) called Git. Git is responsible for everything GitHub-related that happens locally on your computer.

To use Git on the command line, you'll need to download, install, and configure Git on your computer.

If you want to work with Git locally, but don't want to use the command line, you can instead download and install the GitHub Desktop client.

At the heart of GitHub is an open source version control system (VCS) called Git. Git is responsible for everything GitHub-related that happens locally on your computer.

To use Git on the command line, you'll need to download, install, and configure Git on your computer.
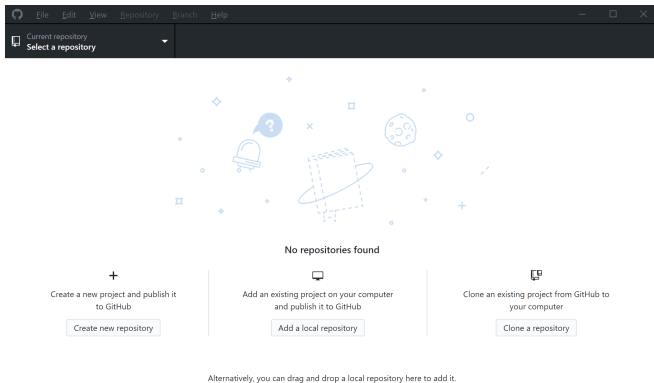
If you want to work with Git locally, but don't want to use the command line, you can instead download and install the GitHub Desktop client.

If you don't need to work with files locally, GitHub lets you complete many Git-related actions directly in the browser, including:

# Git/Github

At the heart of GitHub is an open source version control system (VCS) called Git. Git is responsible for everything GitHub-related that happens locally on your computer.

To use Git on the command line, you'll need to download, install, and configure Git on your computer.

If you want to work with Git locally, but don't want to use the command line, you can instead download and install the GitHub Desktop client.

If you don't need to work with files locally, GitHub lets you complete many Git-related actions directly in the browser, including:

- Creating a repository - Forking a repository - Managing files - Being social
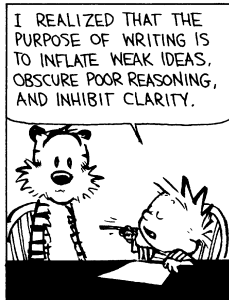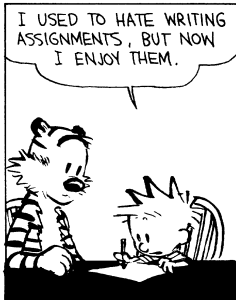
https://desktop.github.com/

What problems does reproducibility solve?

- Transparency
- Data availability
- Software/methods availability
- Improved transfer of knowledge
- ?Validity / correctness of the analysis

# References

Christopher Gandrud (2014). Reproducible Research with R and RStudio, CRC Press.

Victoria Stodden, Friedrich Leisch, Roger Peng (2014). Implementing Reproducible Research. CRC Press.

Yihui Xie. Dynamic Documents with R and knitr, CRC Press

You can start looking at this cheat sheet:
https://www.rstudio.com/wp-content/uploads/2015/02/rmarkdown-cheatsheet.pdf.

For more details on using R Markdown see
http://rmarkdown.rstudio.com.

How to setup Github using web:
https://dannguyen.github.io/github-for-portfolios/lessons/setup-github/.

For more details on using R Markdown see
http://rmarkdown.rstudio.com.

https://sebnemer.github.io/