

Reporting using Dynamic Documents

L. Torgo

ltorgo@fc.up.pt

Faculdade de Ciências / LIAAD-INESC TEC, LA
Universidade do Porto

Jun, 2017



Introduction

Reporting

The Standard Process of Data Analysis

- Import the data into our favorite data analysis tool
- Carry out a set of data analysis steps
- Report the work by building some document (report and / or presentation)
 - Series of *copy+paste* steps from parts of the results of the analysis into some word processing and/or presentation software tool, adding some supporting text
- Frequently all process needs to be repeated / iterated!

Some of the Dangers of this Standard Approach

- Too many manual steps \mapsto great potential for human error
- Too much human effort (time) in the process with many repetitive and boring tasks, like for instance the communication between different software tools
- All process is hardly recordable for future re-use, due to the extensive use of graphical user interfaces
- Small changes on the initial data require full repetition of all process!
- The Analysis and the Reporting are separated and thus great care is required to have both in “sync” avoiding reporting errors
- Very hard to share the work with other teams
- Very hard to re-use the work on similar tasks

List inspired by *Dynamic Documents with R and knitr* by Yihui Xie

Dynamic Documents

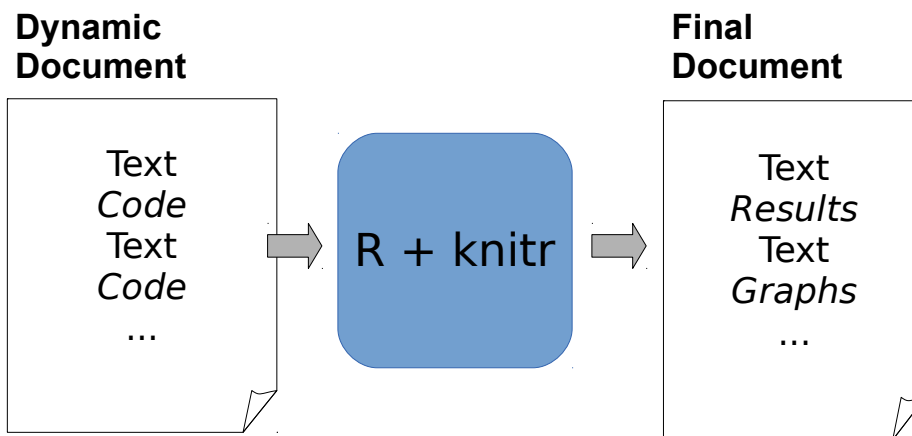
What are Dynamic Documents

- Documents that mix data analysis steps with descriptive text
- Documents that are executable by a computer program to produce the final document
- This final document is produced by a computer program from the initial document created by the user containing data analysis steps and descriptive text

Dynamic Documents solve most of the problems we have described before!

Dynamic Documents

R+knitr - an implementation of the idea



STERN

MS in Business Analytics

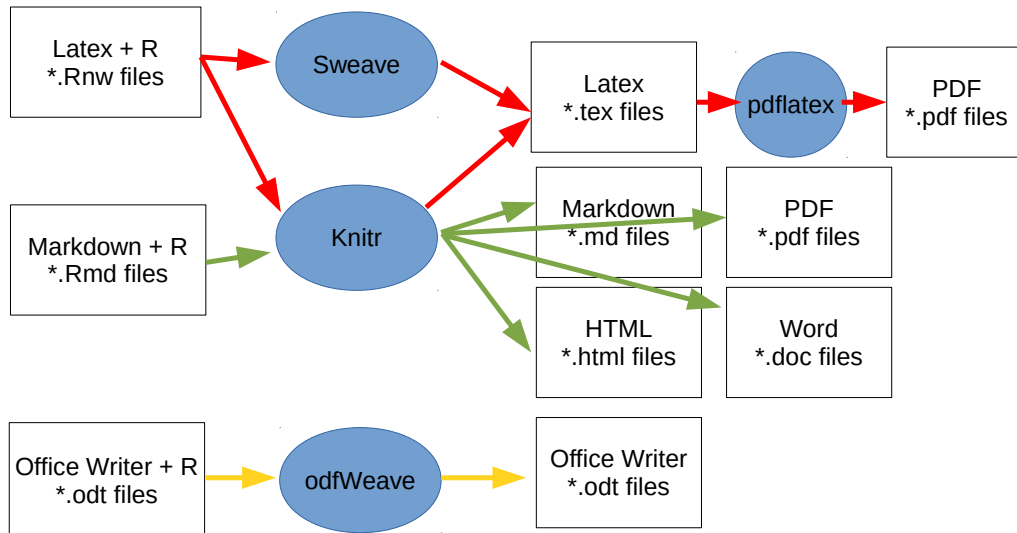
Dynamic Documents

- The idea is related with the concept of *literate programming*
- It is implemented in 3 main steps:
 - 1 Inspect the dynamic report and separate the data analysis code from the descriptive text
 - 2 Execute the code and store the results of each code chunk
 - 3 Produce the final document replacing the data analysis code in the original document by the results
- NOTE: all process is carried out without human intervention!

Knuth, Donald E. (1984). "Literate Programming". The Computer Journal (British Computer Society) 27 (2): 97–111.

Dynamic Documents

Different implementations of the concept



STERN

MS in Business Analytics

R+knitr

- Knitr is a platform that allows:
 - Using different document formats
 - \LaTeX , HTML and Markdown
 - Using different programming languages to implement the data analysis steps
 - R, Python, awk, C++ and shell scripts

rmarkdown

- **rmarkdown** is an R package that allows integrating R with Markdown using **knitr** and **pandoc**
- **Pandoc** is a kind of swiss-army knife of document format conversion
- **R markdown** is a literate programming approach that allows embedding R code into markdown documents
- **Markdown** is a very simple markup language that was designed to easily produce internet content

A (Very) Brief Introduction to Markdown

- Very simple language for creating Internet content
 - Much simpler than the native language - HTML
 - A tool with the same name can be used to convert markdown into HTML
- The format is a simple text file and you do not need any special software tool to create Markdown files

Character formatting in Markdown

- Formatting is carried out with the help of small annotation tags

- Examples:

`**Note**` or `__Note__`

translated into **Note** (i.e. boldface)

Sections and subsections

- An hash (#) character before a text line indicates a first level section heading
- If instead you use two or three hashes we get second and third level sections

Illustrations on the use of Markdow in Dynamic Reports

Introduction

R comes with a series of data sets. We may use the function `data` to load them, as show in the following example:

```
# Illustrations on the use of Markdow in Dynamic Reports
```

```
## Introduction
```

```
R comes with a series of data sets. We may use the function
**data** to load them, as show in the following example:
```

Lists of Items

OM prototype - Admin User

+ Stage 1

- Define topics of interest
- Define web sources
 - * create crawlers (**Potentially Challenging)
 - * collect data

+ Stage 2

- Obtain tagged data (**Potentially Challenging)
- Obtain, evaluate and select models

OM prototype - Admin User

- Stage 1
 - Define topics of interest
 - Define web sources
 - create crawlers (Potentially Challenging)
 - collect data
- Stage 2
 - Obtain tagged data (Potentially Challenging)
 - Obtain, evaluate and select models



Embed images and links

Monitoring and Forecasting Water Quality Parameters

- Collaboration with [Águas do Douro e Paiva, SA] (<http://addp.pt/pt/home.php>)
- FCT project [MORWAQ] (<http://liaad.inescporto.pt/modys/projects/morwaq>)
- Some of the main results:
 - Software prototype for monitoring and forecasting water quality parameters
 - Several publications
 - KDD'2011, ECAI'2010

![ADdP Network] (addpNet.png "The ADdP Network")

Monitoring and Forecasting Water Quality Parameters

- Collaboration with [Águas do Douro e Paiva, SA](#)
- FCT project [MORWAQ](#)
- Some of the main results:
 - Software prototype for monitoring and forecasting water quality parameters
 - Several publications
 - KDD'2011, ECAI'2010



R Markdown

Embedding R code into documents

The Distribution of Sepal Length

An histogram of the variable
can be obtained with:

```
```{r eval=FALSE}
hist(iris$Sepal.Length)
```
```

```
```{r echo=FALSE}
data(iris)
hist(iris$Sepal.Length)
```
```

Chunk+Insert Chunk in RStudio

The Distribution of Sepal Length

An histogram of the variable can be obtained with:

```
hist(iris$Sepal.Length)
```



Code Chunks in R Markdown

- Code chunks are parts of a document that should be executed by R
- In R Markdown they are indicated as follows:


```
```{r}
...
```
```
- Everything between these two lines will be executed by R
- The result of the execution will be part of the final document
- Inline code can be used like this: ``r` 2+2``

Code Chunk Options

- The way code chunks are interpreted by Knitr is controllable through a series of chunk options
- Without any option code chunks are:
 - 1 Executed
 - 2 A code block is added to the final document
 - 3 A results block is added to the final block
- Chunk options allow to adjust these defaults
- This can be done on each individual chunk or globally

Some common chunk options

eval (TRUE) - allows to control whether the chunk code is to be executed or not by R. In the following example the code is inserted into the final document but it will not be executed by R

```
```{r eval=FALSE}
hist(iris$Sepal.Length)
```

%$
```

echo (TRUE) - allows to hide the R code from the final document (with the value FALSE), with only the results of the execution being included in that document (e.g. a figure)

- Note that each chunk may include several options separated by commas
- Many more options exist! See an exhaustive list with explanations at <http://yihui.name/knitr/options>

Generating the Final Document

- Function `render()` from package **rmarkdown** receives as input a markdown file and produces the final document using R to execute the code chunks:

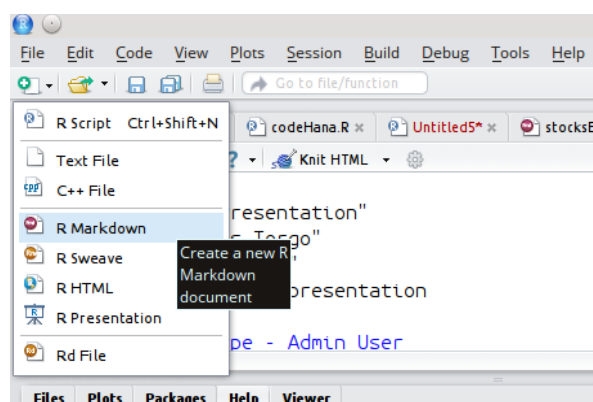
```
library(rmarkdown)
render("initialDoc.Rmd")
```

- The type of output document that is generated can be controlled through meta-data information that is included in the initial document,

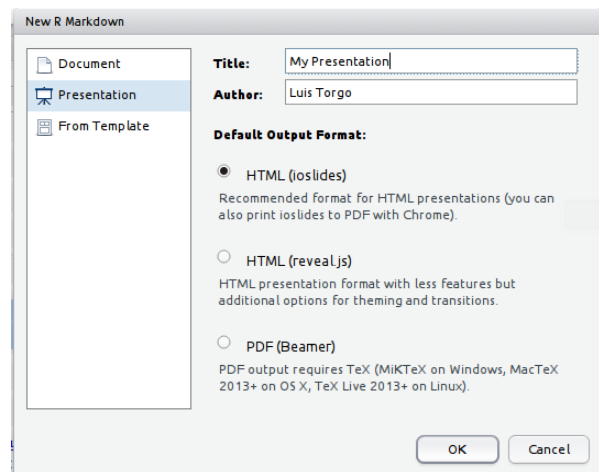
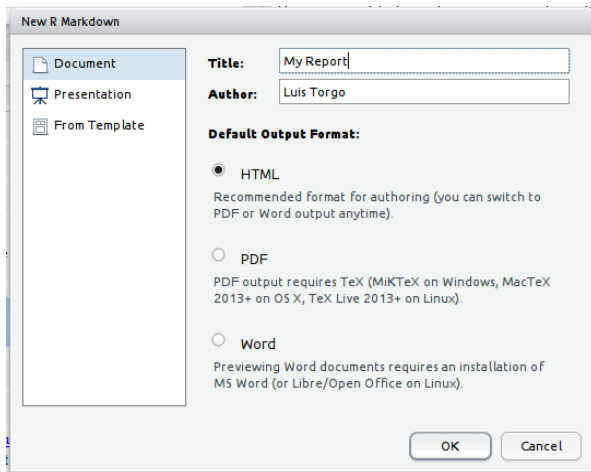
```
---
title: "My First Dynamic Report"
output:
  pdf_document:
    toc: true
    highlight: zenburn
---
```

Using RStudio interface

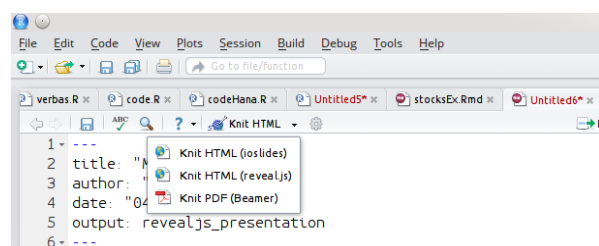
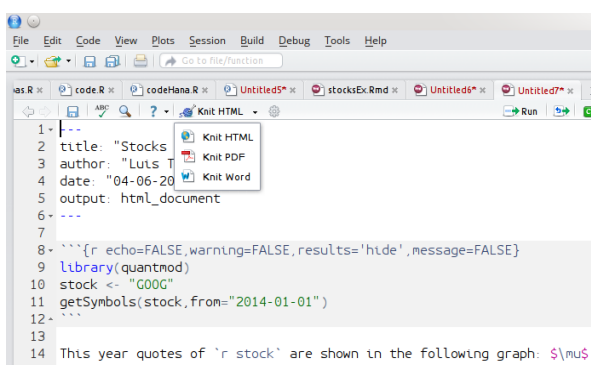
Recent versions of RStudio facilitate the creation of dynamic reports and presentations using R Markdown



Selecting the Type of Document Using RStudio interface



Generating the Output Document Using RStudio interface



Learning more

Much more information is available at
<http://rmarkdown.rstudio.com/>

R Notebooks

R Notebooks

- An R Notebook are a new feature of RStudio
- They are a type of R markdown document.
- Its main characteristics are that code chunks that can be executed independently and interactively, with output visible immediately below the chunk.
- Any R markdown document can be used as a R notebook in the recent versions of RStudio
- Any R Notebook can be rendered to some of the output formats available to R Markdown documents
- R Notebooks are a nice way of interaction with R and at the same time documenting your activity
- As any R markdown document, R Notebooks are also easily shareable with your collaborators

An example of an R Notebook

The screenshot shows an R Notebook interface. The top section contains a YAML header and introductory text. The bottom section contains an R code chunk that has been executed, resulting in a scatter plot.

```

1 ---
2 title: "R Notebook"
3 output: html_notebook
4 ---
5
6 This is an [R Markdown](http://rmarkdown.rstudio.com) Notebook. When you execute code within the
7 notebook, the results appear beneath the code.
8
9 Try executing this chunk by clicking the *Run* button within the chunk or by placing your cursor
10 inside it and pressing *Ctrl+Shift+Enter*.
11
12 ```{r}
13 plot(cars)
14 ```

```

The scatter plot displays the relationship between speed (x-axis, ranging from 5 to 25) and distance (y-axis, ranging from 0 to 120) for the 'cars' dataset. The plot shows a positive correlation, with distance increasing as speed increases.

Below the plot, the following text is visible:

```

13
14 Add a new chunk by clicking the *Insert Chunk* button on the toolbar or by pressing *Ctrl+Alt+I*.
15
16 When you save the notebook, an HTML file containing the code and output will be saved alongside it
17 (click the *Preview* button or press *Ctrl+Shift+K* to preview the HTML file).

```

Creating R Notebooks

- Using RStudio menus: *File -> New File-> R Notebook*
- Using `html_document` in your YAML document header:

```

---
title: "Day 1 Classes Log"
output: html_notebook
---

```

Saving and Sharing R Notebooks

- When you save the R Notebook two files are created:
 - The standard `.Rmd` (R markdown) file
 - An html (with extension `.nb.html`) that can be opened on any browser to visualize the notebook
- You may share your notebook through both files
- Actually, as the `.nb.html` file contains the `.Rmd` file, you may share only this file and when your colleagues open it in RStudio it will extract automatically the `Rmd` file and open it in the editor!