

# Reproducible Research with R

Arnaud Legrand

COMPAS tutorials, Neuchâtel, April 2014

## 1 Reproducible Research

- Looks familiar ?
- Many Different Alternatives

## 2 R

- General Introduction
- Reproducible Documents: knitR
- Introduction to R
- Needful Packages by Hadley Wickam

## 1 Reproducible Research

- Looks familiar ?
- Many Different Alternatives

## 2 R

- General Introduction
- Reproducible Documents: knitR
- Introduction to R
- Needful Packages by Hadley Wickam

This may be an interesting contribution but:

- This **average value** must hide something.
- As usual, there is no **confidence interval**, I wonder about the variability and whether the difference is **significant** or not.
- That can't be true, I'm sure they **removed some points**.
- Why is this graph in **logscale** ? How would it look like otherwise ?
- The authors decided to show only a **subset of the data**. I wonder what the rest looks like.
- There is no label/legend/... What is the **meaning of this graph** ? If only I could access the generation script.

- I thought I used the same parameters but I'm getting different results!
- The new student wants to compare with the method I proposed last year.
- The damn reviewer asked for a major revision and wants me to change this figure. :(
- Which code and which data set did I use to generate this figure?
- It worked yesterday!
- Why did I do that?

# My Feeling

Computer scientists have an incredibly **poor training in probabilities and statistics**

Why should we ? Computer are **deterministic** machines after all, right? ;)

Eight years ago, I've started realizing how **lame** the articles I reviewed (as well as those I wrote) were in term of experimental methodology.

- Yeah, I know, your method/algorithm is better than the others as demonstrated by the figures
- Not enough information to **discriminate real effects from noise**
- Little information about the **workload**
- Would the “conclusion” still hold with a slightly different workload?
- I'm tired of awful combination of tools (perl, gnuplot, sql, ...) and **bad methodology**

# Current practice in CS

Computer scientists tend to either:

- vary **one factor at a time**, use a very fine sampling of the parameter range,
- **run millions of experiments** for a week varying a lot of parameters and then try to get something of it. Most of the time, they (1) don't know how to analyze the results (2) realize something went wrong. . .

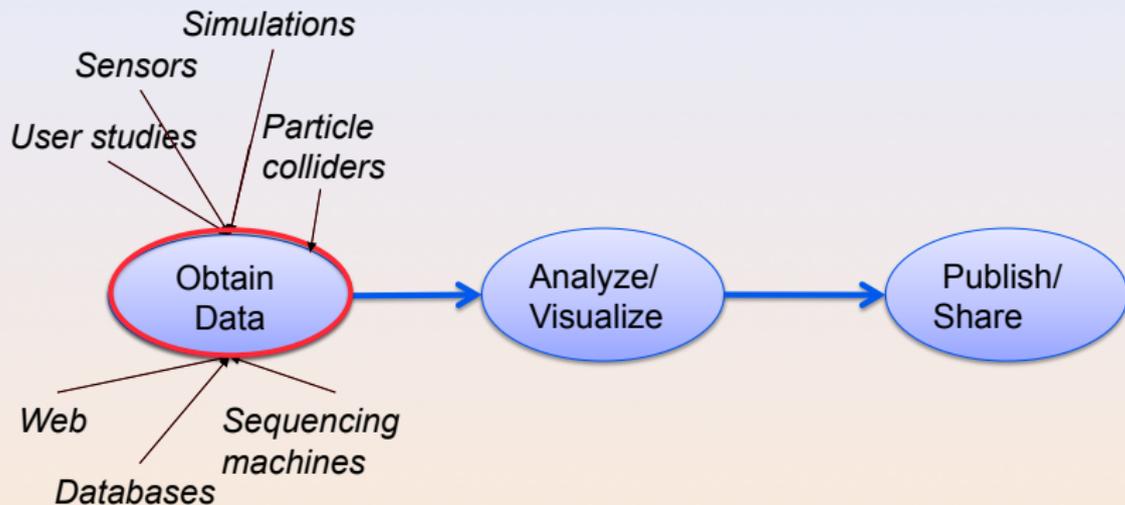
Interestingly, most other scientists do **the exact opposite**.

These two flaws come from poor training and from the fact that C.S. experiments are **almost** free and very fast to conduct.

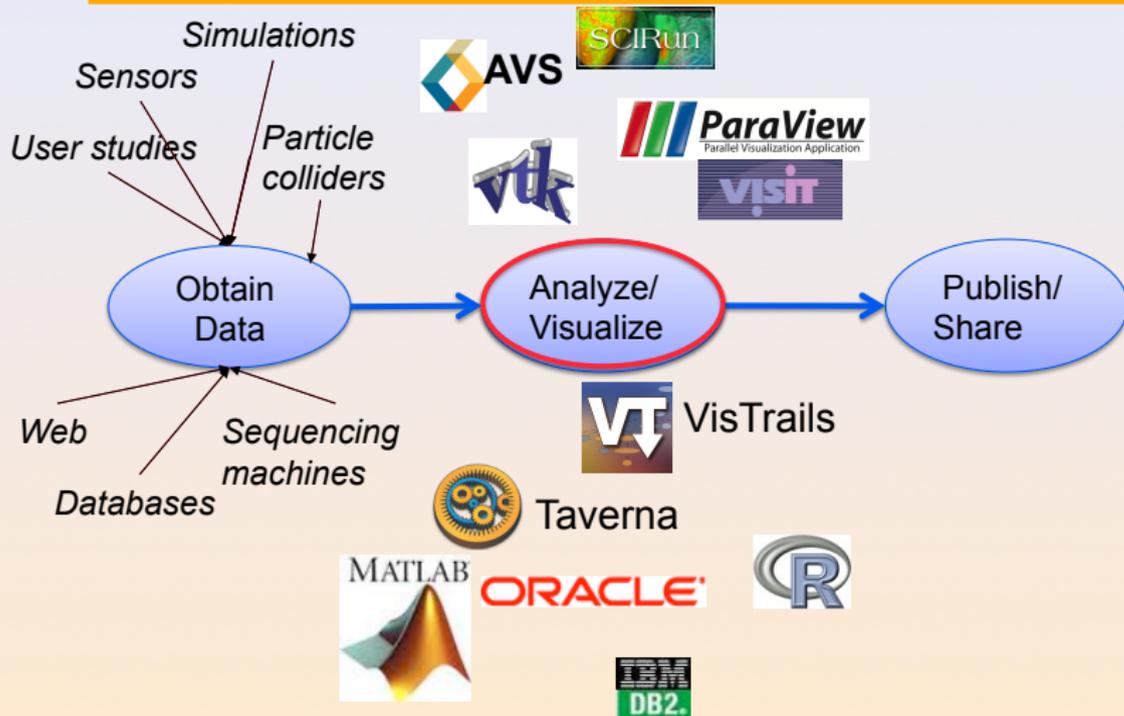
- Most strategies of experimentation have been designed to **provide sound answers despite** all the **randomness and uncontrollable factors**;
- **Maximize the amount of information** provided by a given set of experiments;
- **Reduce** as much as possible **the number of experiments** to perform to answer a given question under a given level of confidence.

Takes a few lectures on **Design of Experiments** to improve but anyone can start by reading **Jain's book on The Art of Computer Systems Performance Analysis**

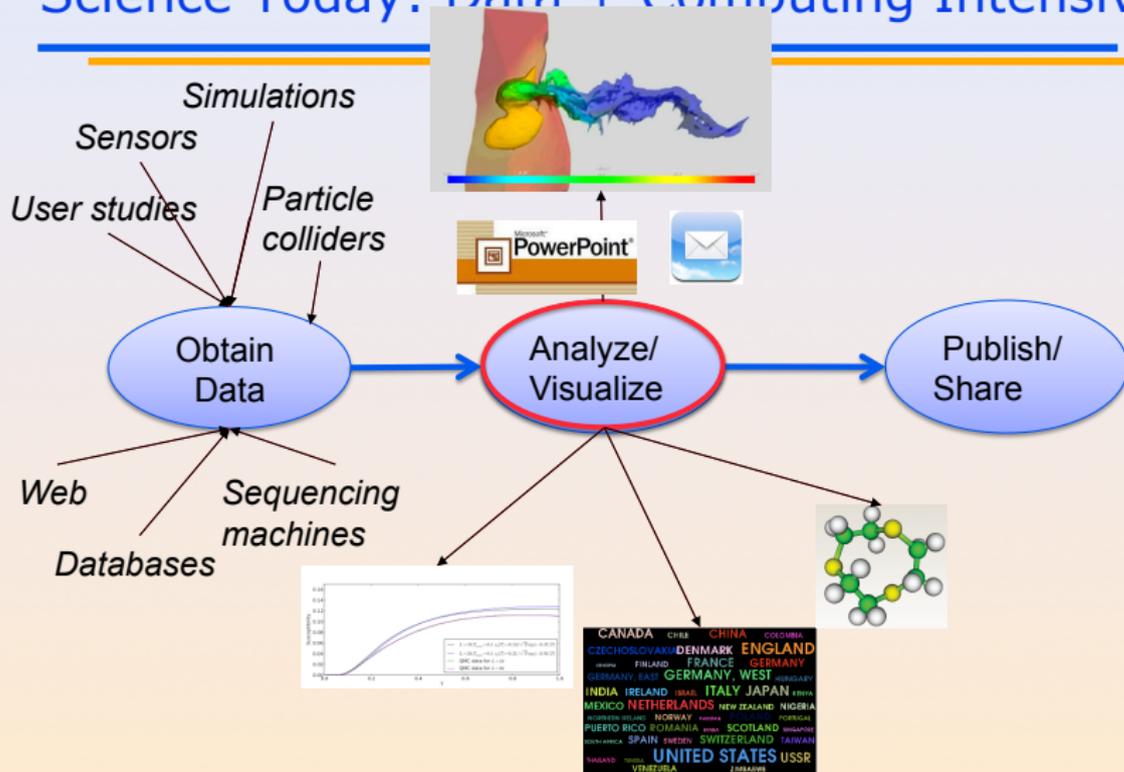
# Science Today: Data Intensive



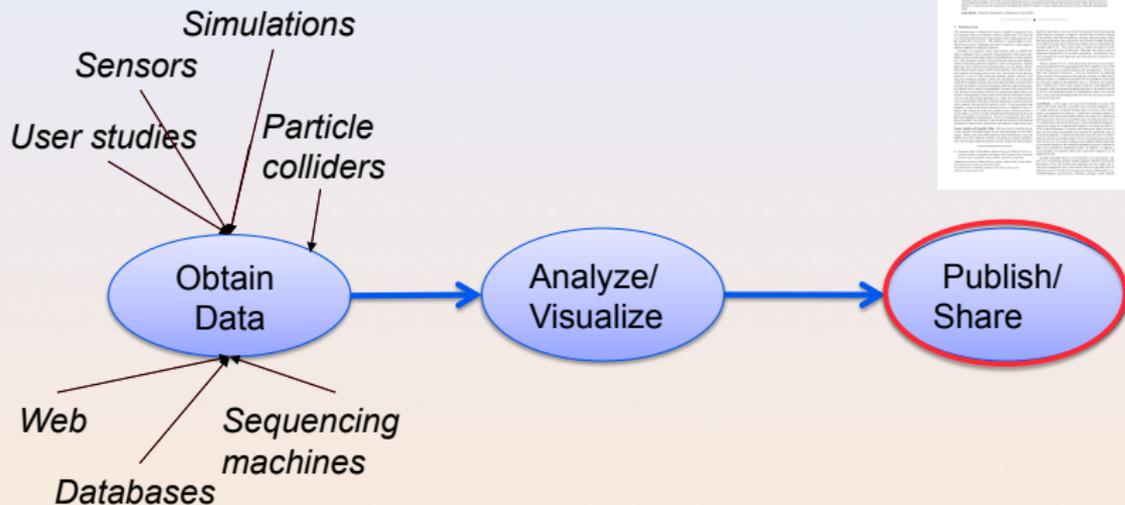
# Science Today: Data + Computing Intensive



# Science Today: Data + Computing Intensive



# Science Today: Data + Computing Inte



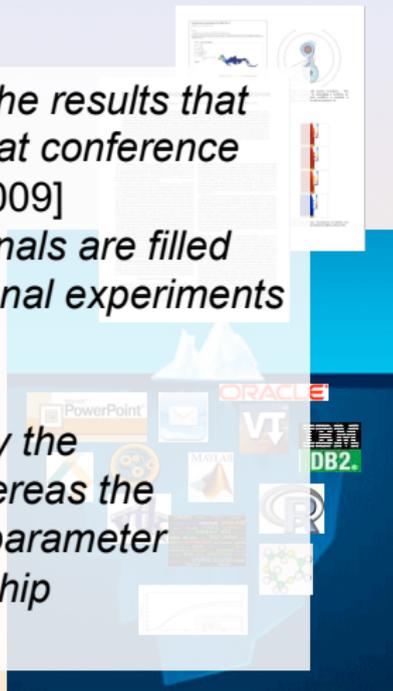
# Science Today: Incomplete Publications

- ◆ Publications are just the tip of the iceberg
  - Scientific record is incomplete---to large to fit in a paper
  - Large volumes of data
  - Complex processes
- ◆ Can't (easily) reproduce results

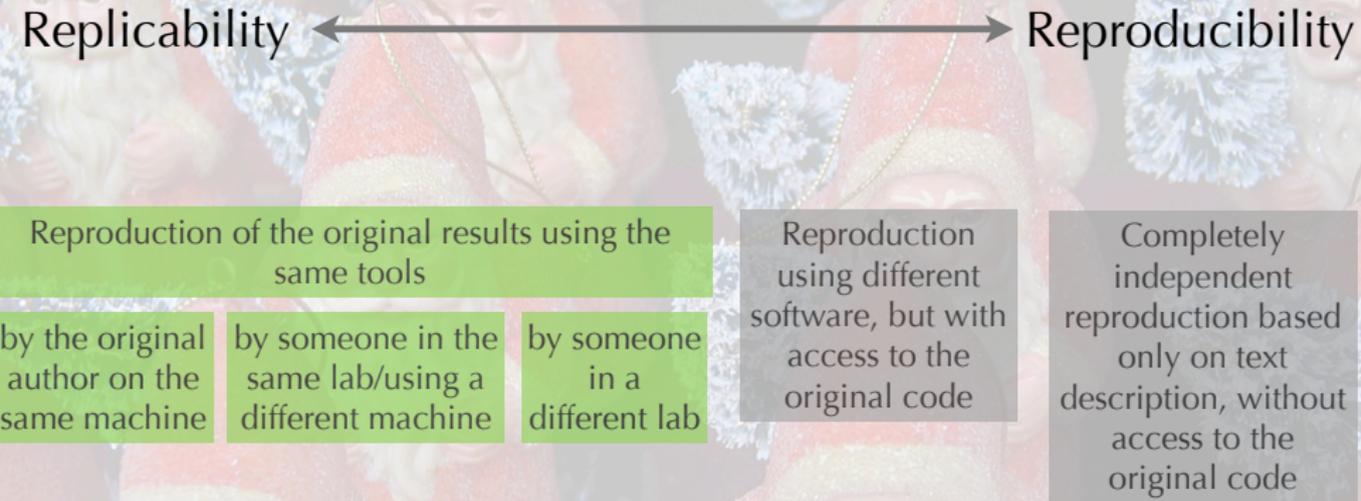


# Science Today: Incomplete Publications

- ◆ Publications are just the tip of the iceberg
  - *“It’s impossible to verify most of the results that computational scientists present at conference and in papers.”* [Donoho et al., 2009]
  - *“Scientific and mathematical journals are filled with pretty pictures of computational experiments that the reader has no hope of repeating.”* [LeVeque, 2009]
  - *“Published documents are merely the advertisement of scholarship whereas the computer programs, input data, parameter values, etc. embody the scholarship itself.”* [Schwab et al., 2007]



# Reproducibility: What Are We Talking About ?



# A Difficult Trade-off

## Automatically keeping track of everything

- the code that was run (source code, libraries, compilation procedure)
- processor architecture, OS, machine, date, . . .

VM-based solutions

## Ensuring others can redo/understand what you did

- Why did I run this?
- Does it still work when I change this piece of code for this one?

Laboratory notebook and recipes

# Reproducible Research: the New Buzzword ?

## H2020-EINFRA-2014-2015

*A key element will be capacity building to link literature and data in order to enable a more transparent evaluation of research and **reproducibility** of results.*

## More and more workshops

- Workshop on Duplicating, Deconstructing and Debunking (WDDD) (2014 edition)
- **Reproducible Research: Tools and Strategies for Scientific Computing** (2011)
- Working towards Sustainable Software for Science: Practice and Experiences (2013)
- **REPPAR'14: 1st International Workshop on Reproducibility in Parallel Computing**
- Reproducibility@XSEDE: An XSEDE14 Workshop
- Reproduce/HPCA 2014
- TRUST 2014

Should be seen as opportunities to share experience.

## 1 Reproducible Research

- Looks familiar ?
- Many Different Alternatives

## 2 R

- General Introduction
- Reproducible Documents: knitR
- Introduction to R
- Needful Packages by Hadley Wickam

## Our Approach: An Infrastructure to Support Provenance-Rich Papers [Koop et al., ICCS 2011]

---

- ◆ Tools for *authors* to create reproducible papers
  - Specifications that encode the computational processes
  - Package the results *Support different approaches*
  - Link from publications
- ◆ Tools for testers to repeat and validate results
  - Explore different parameters, data sets, algorithms
- ◆ Interfaces for searching, comparing and analyzing experiments and results
  - Can we discover better approaches to a given problem?
  - Or discover relationships among workflows and the problems?
  - How to describe experiments?

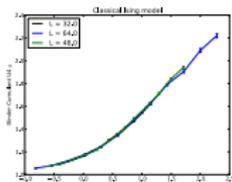
## An Provenance-Rich Paper: ALPS2.0

arXiv:1101.2646v4 [cond-mat.str-el] 23 May 2011

The ALPS project release 2.0:  
Open source software for strongly correlated  
systems

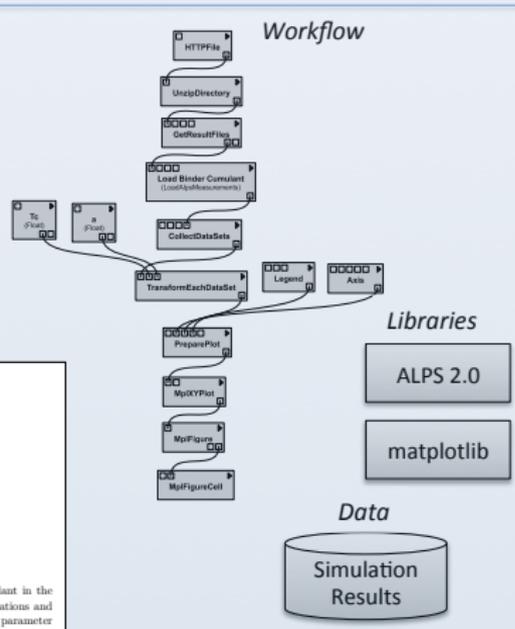
B. Bauer<sup>1</sup> L. D. Carr<sup>2</sup> H.G. Evertz<sup>3</sup> A. Feiguin<sup>4</sup> J. Freire<sup>5</sup>  
S. Fuchs<sup>6</sup> L. Gamper<sup>7</sup> J. Gukelberger<sup>8</sup> E. Gull<sup>9</sup> S. Guertler<sup>4</sup>  
A. Hehn<sup>10</sup> R. Igarashi<sup>11,10</sup> S.V. Isakov<sup>1</sup> D. Koop<sup>1</sup> P.N. Ma<sup>1</sup>  
P. Mates<sup>12</sup> H. Matsuo<sup>13</sup> O. Parcollet<sup>12</sup> G. Pawłowski<sup>13</sup>  
J.D. Picon<sup>14</sup> L. Pollet<sup>15</sup> E. Santos<sup>6</sup> V.W. Scarola<sup>16</sup>  
U. Schollwöck<sup>17</sup> C. Silva<sup>8</sup> B. Surer<sup>8</sup> S. Todo<sup>18,11</sup> S. Trebst<sup>18</sup>  
M. Troyer<sup>1</sup> M. L. Wall<sup>1</sup> P. Werner<sup>8</sup> S. Wessel<sup>19,20</sup>

- <sup>1</sup>Theoretische Physik, ETH Zurich, 8005 Zurich, Switzerland
- <sup>2</sup>Department of Physics, Colorado School of Mines, Golden, CO 80401, USA
- <sup>3</sup>Institut für Theoretische Physik, Technische Universität Graz, A-8010 Graz, Austria
- <sup>4</sup>Department of Physics and Astronomy, University of Wyoming, Laramie, Wyoming 82071, USA
- <sup>5</sup>Scientific Computing and Imaging Institute, University of Utah, Salt Lake City, Utah 84112, USA
- <sup>6</sup>Institut für Theoretische Physik, Georg-August-Universität Göttingen, Göttingen, Germany
- <sup>7</sup>Columbia University, New York, NY 10027, USA
- <sup>8</sup>Bethe Center for Theoretical Physics, Universität Bonn, Nussallee 12, 53115 Bonn, Germany



1 Correspond

**Figure 3.** In this example we show a data collapse of the Binder Cumulant in the classical Ising model. The data has been produced by remotely run simulations and the critical exponent has been obtained with the help of the VisTrails parameter exploration functionality.



## Chronicling computations in real-time

VCR computation platform Plugin = Computation recorder

### Regular program code

```
figure1 = plot(x)
save(figure1, 'figure1.eps')
```

```
> file /home/figure1.eps saved
>
```

## Chronicling computations in real-time

VCR computation platform Plugin = Computation recorder

Program code with VCR plugin

```
repository vcr.nature.com  
verifiable figure1 = plot(x)
```

```
> vcr.nature.com approved:  
> access figure1 at https://vcr.nature.com/ffaaffb148d7
```

## Word-processor plugin App

### LaTeX source

```
\includegraphics{figure1.eps}
```

### LaTeX source with VCR package

```
\includeresult{vcr.thelancet.com/ffaaffb148d7}
```

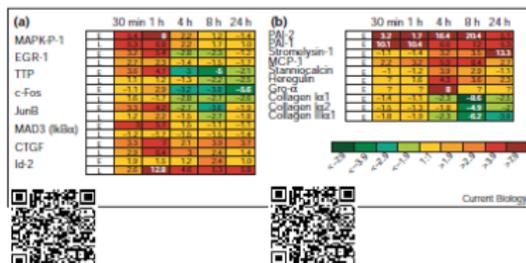
Permanently bind printed graphics to underlying result content

# VCR: A Universal Identifier for Computational Results

Research Paper Analysis of replicative senescence Shelton et al. 943

Figure 3

Time course of serum stimulation. (a) Early passage (E: PD30) or late passage (L: PD89) BJ cultures were held in 0.5% serum for 2 days, then stimulated with 10% FBS. RNA levels from cultures at the indicated time points (Cy5 channel) were compared with the uninduced starting culture (Cy3 channel). Positive values indicate higher expression in induced cells; negative values indicate lower expression in induced cells. Question marks indicate that there was insufficient signal for detection. A complete listing of serum-responsive genes from this analysis is provided in Supplementary material. (b) The serum-responsiveness of select senescence-regulated genes in early passage (PD30) BJ fibroblasts.



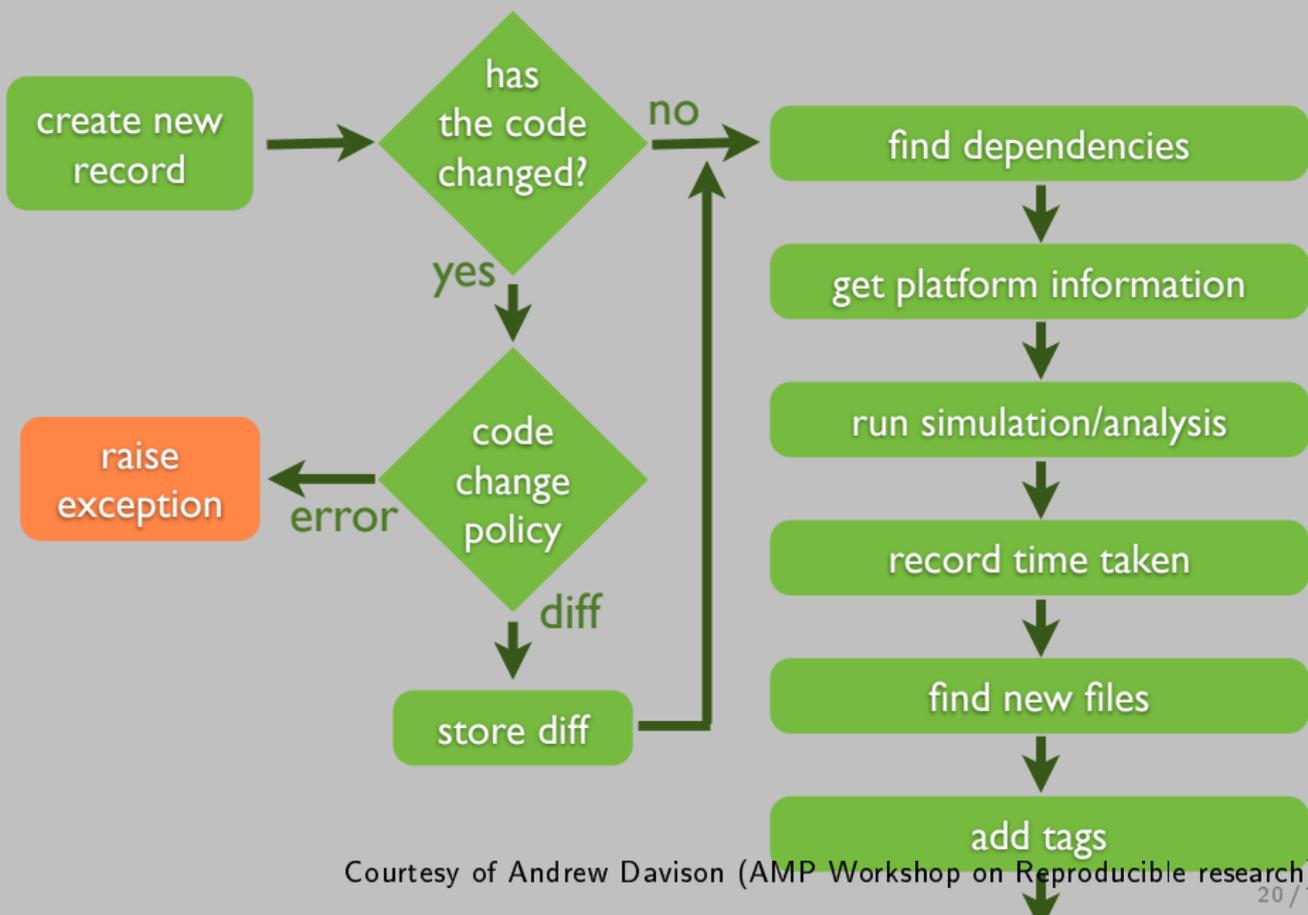
senescence response appears to overlap substantially with gene expression patterns observed in activated fibroblasts during wound healing [24–26]. MCP-1, Gro- $\alpha$ , IL-1 $\beta$  and IL-15 are strong effectors of macrophage and neutrophil recruitment and activation [27,28]. The upregulation of Toll (Tlr-4) in senescent fibroblasts confirms the overall immune response behavior at senescence. Tlr-4 is an IL-1 receptor homolog and is implicated in the activation of the gene regulatory protein NF- $\kappa$ B, a function proposed to be part of the innate immune response [29]. The induction of IL-15 at senescence is also consistent with an innate immune response, as IL-15 can be induced by NF- $\kappa$ B-dependent transcription [30] and also participates in inflammatory disease processes [28].

Deficiencies in the response of senescent cells to serum stimulation have been reported, and include an inability to induce the expression of *c-fos* mRNA [31] and markers of late G1 and S phase [32]. In response to serum, expression of inflammatory chemokines, matrix-degrading proteases and their modulators is induced in early-passage dermal fibroblasts, and expression of matrix collagens is reduced. This transient burst of activity may represent the natural response of the cells to wound repair [24]. Id-2 transcripts were hyper-induced in serum-stimulated senescent

states overlap substantially with those in telomere-induced senescence (W.F., D.N.S., R. Allsopp, S. Lowe, and G. Ferbeyre, unpublished observations) and thus are likely to use many of the same activation processes.

The pattern of gene expression at senescence varies substantially in different cell types. Although the expression of matrix and structural proteins, such as the collagens, keratins and auxiliary factors, is repressed in RPE cells, inflammatory regulators are not induced, in contrast to dermal fibroblasts. Physiologically, this would make sense, as an acute inflammatory response in a tissue critical for normal vision would be likely to have deleterious consequences. However, as the RPE layer has a central role in the deposition and maintenance of extracellular matrix in the retina, decrements in the ability of senescent RPE cells to maintain appropriate expression patterns, as evidenced by decreased expression of collagens, keratins, aggrecan, transglutaminase and so on, would be predicted to have adverse effects on retinal architecture. Dysfunction of the RPE cell layer is considered to be a substantial factor in the development of age-related macular degeneration [36].

# Sumatra: a lab notebook



Courtesy of Andrew Davison (AMP Workshop on Reproducible research)

```
$ smt comment 20110713-174949 "Eureka! Nobel prize  
here we come."
```

```
$ smt tag "Figure 6"
```

# Sumatra: a lab notebook

Sumatra: TestProject: List of records

http://127.0.0.1:8002/ Google

### TestProject: List of records

Delete Include data	Label	Reason	Outcome	Duration	Processes	Simulator		Script			Date	Time	Tags
						Name	Version	Repository	Main file	Version			
<input type="checkbox"/>	<a href="#">20100709-154255</a>		'Eureka! Nobel prize here we come.'	0.59 s		Python	2.5.2	/Users/andrew/tmp/SumatraTest	main.py	396c2020ca50	09/07/2010	15:42:55	
<input type="checkbox"/>	<a href="#">20100709-154309</a>			0.59 s		Python	2.5.2	/Users/andrew/tmp/SumatraTest	main.py	396c2020ca50	09/07/2010	15:43:09	
<input type="checkbox"/>	<a href="#">hagging</a>	'determine whether the gourd is worth 3 or 4 shekels'	'apparently, it is worth NaN shekels.'	0.59 s		Python	2.5.2	/Users/andrew/tmp/SumatraTest	main.py	396c2020ca50	09/07/2010	15:43:20	fooba
<input type="checkbox"/>	<a href="#">20100709-154338</a>	'test effect of a smaller time constant'		0.59 s		Python	2.5.2	/Users/andrew/tmp/SumatraTest	main.py	396c2020ca50	09/07/2010	15:43:38	
<input type="checkbox"/>	<a href="#">hagging_repeat</a>	Repeat experiment hagging	The new record exactly matches the original.	0.58 s		Python	2.5.2	/Users/andrew/tmp/SumatraTest	main.py	396c2020ca50	09/07/2010	15:43:47	

Courtesy of Andrew Davison (AMP Workshop on Reproducible research)

# New Tools for Computational Reproducibility

- Dissemination Platforms:

[ResearchCompendia.org](http://ResearchCompendia.org)

[IPOL](http://IPOL)

[Madagascar](http://Madagascar)

[MLOSS.org](http://MLOSS.org)

[thedatahub.org](http://thedatahub.org)

[nanoHUB.org](http://nanoHUB.org)

[Open Science Framework](http://Open Science Framework)

[The DataVerse Network](http://The DataVerse Network)

[RunMyCode.org](http://RunMyCode.org)

- Workflow Tracking and Research Environments:

[VisTrails](http://VisTrails)

[Kepler](http://Kepler)

[CDE](http://CDE)

[Galaxy](http://Galaxy)

[GenePattern](http://GenePattern)

[Synapse](http://Synapse)

[Sumatra](http://Sumatra)

[Taverna](http://Taverna)

[Pegasus](http://Pegasus)

- Embedded Publishing:

Courtesy of Victoria Stodden (UC Davis, Feb 13, 2014)

[Verifiable Computational Research](http://Verifiable Computational Research)

[Sweave](http://Sweave)

[knitr](http://knitr)

[Collage Authoring Environment](http://Collage Authoring Environment)

[SHARE](http://SHARE)

And also: Figshare, ActivePapers, Elsevier executable paper, ...

# Literate programming

Donald Knuth: explanation of the program logic in a natural language interspersed with snippets of macros and traditional source code.

I'm way too stupid to program this way but that's exactly what we need for writing a reproducible article/analysis!

## Org-mode (requires emacs)

My favorite tool.

- plain text, very smooth, works both for html, pdf, ...
- allows to combine all my favorite languages

## lpython notebook

If you are a python user, go for it! Web app, easy to use/setup...

## KnitR (a.k.a. Sweave)

For non-emacs users and as a first step toward reproducible papers:

- Click and play with a modern IDE

## 1 Reproducible Research

- Looks familiar ?
- Many Different Alternatives

## 2 R

- **General Introduction**
- Reproducible Documents: knitR
- Introduction to R
- Needful Packages by Hadley Wickam

# Why R?

R is a great language for data analysis and statistics

- Open-source and multi-platform
- Very expressive with high-level constructs
- Excellent graphics
- Widely used in academia and business
- Very active community
  - Documentation, FAQ on <http://stackoverflow.com/questions/tagged/r>
- Great integration with other tools

# Why is R a pain for computer scientists ?

- R is **not** really a **programming** language
- Documentation is for statisticians
- Default plots are cumbersome (meaningful)
- Summaries are cryptic (precise)
- **Steep learning curve** even for us, computer scientists whereas we generally switch seamlessly from a language to another! That's frustrating! ;)

# Do's and don't's

~~R is high level, I'll do everything myself~~

- CTAN comprises 4,334 T<sub>E</sub>X, L<sup>A</sup>T<sub>E</sub>X, and related packages and tools. Most of you do not use plain T<sub>E</sub>X.
- Currently, the CRAN package repository features 4,030 available packages.
- How do you know which one to use ??? Many of them are highly exotic (not to say useless to you).

I learnt with <http://www.r-bloggers.com/>

- Lots of introductions but not necessarily what you're looking for so I'll give you a short tour.  
You should quickly realize though that you need proper training in statistics and data analysis if you do not want tell nonsense.
- Again, you should read Jain's book on *The Art of Computer Systems Performance Analysis*
- You may want to follow online courses:
  - <https://www.coursera.org/course/compdata>
  - <https://www.coursera.org/course/repdata>

# Install and run R on debian

```
1 apt-cache search r
```

Err, that's not very useful :) It's the same when searching on google but once the filter bubble is set up, it gets better...

```
1 sudo apt-get install r-base
```

```
1 R
```

```
1 R version 3.0.2 (2013-09-25) -- "Frisbee Sailing"
2 Copyright (C) 2013 The R Foundation for Statistical Computing
3 Platform: x86_64-pc-linux-gnu (64-bit)
4
5 R is free software and comes with ABSOLUTELY NO WARRANTY.
6 You are welcome to redistribute it under certain conditions.
7 Type 'license()' or 'licence()' for distribution details.
8
9 R is a collaborative project with many contributors.
10 Type 'contributors()' for more information and
11 'citation()' on how to cite R or R packages in publications.
12
13 Type 'demo()' for some demos, 'help()' for on-line help, or
14 'help.start()' for an HTML browser interface to help.
15 Type 'q()' to quit R.
16 >
```

## Install a few cool packages

R has its own package management mechanism so just run R and type the following commands:

- `ddply`, `reshape` and `ggplot2` by Hadley Wickham (<http://had.co.nz/>)

```
1 install.packages("plyr")
2 install.packages("reshape")
3 install.packages("ggplot2")
```

- `knitr` by (Yihui Xie) <http://yihui.name/knitr/>

```
1 install.packages("knitr")
```

Using R interactively is nice but quickly becomes painful so at some point, you'll want an IDE.

Emacs is great but you'll need *Emacs Speaks Statistics*

```
1 sudo apt-get install ess
```

In this tutorial, we will use **rstudio** (<https://www.rstudio.com/>).

## 1 Reproducible Research

- Looks familiar ?
- Many Different Alternatives

## 2 R

- General Introduction
- **Reproducible Documents: knitR**
- Introduction to R
- Needful Packages by Hadley Wickam

# Rstudio screenshot

The screenshot displays the RStudio interface with three main panels:

- Source Editor:** Contains R code for generating a random variable and plotting it. The code includes comments and R-specific syntax like `knitr` and `knitr::`.
- Environment/History:** Shows the data frame `df` with 10 observations and 2 variables. The variables are `x` (integer[10]) and `y` (numeric[10]).
- Console:** Shows the execution of the R code, including the output of `set.seed(1234)`, `library(ggplot2)`, `library(lattice)`, and the resulting data frame `df` with 10 rows of `x` and `y` values. The final command `plot(x)` is being executed.

The console output for `df` is as follows:

```
> df
  x y
1 1 1.31
2 2 2.31
3 3 3.36
4 4 3.27
5 5 5.04
6 6 6.11
7 7 8.43
8 8 8.98
9 9 8.38
10 10 9.27
```

The plot shows a scatter plot of `x` versus `y`. The x-axis is labeled 'Index' and ranges from 0 to 10. The y-axis is labeled 'x' and ranges from 0 to 10. The data points are approximately:

Index	x
1	1.31
2	2.31
3	3.36
4	3.27
5	5.04
6	6.11
7	8.43
8	8.98
9	8.38
10	9.27

# Reproducible analysis in Markdown + R

- Create a new **R Markdown** document (Rmd) in rstudio
- R chunks are interspersed with “`{r}`” and “”
- Inline R code: `'r sin(2+2)'`
- You can **knit** the document and share it via **rpubs**
- R chunks can be sent to the top-level with **Alt-Ctrl-c**
- I usually work mostly with the current environment and only knit in the end
- Other engines can be used (use rstudio **completion**)

```
1 ‘‘‘{r engine='sh'}
2 ls /tmp/
3 ‘‘‘
```

- Makes **reproducible analysis as simple as one click**
- Great tool for quick analysis for self and colleagues, homeworks, ...

- Create a new **R Sweave** document (Rnw) in rstudio
- R chunks are interspersed with `<<>=` and `@`
- You can **knit** the document to produce a pdf
- You'll probably quickly want to **change default behavior** (activate the cache, hide code, ...). In the preamble:

```
1 <<echo=FALSE>>=  
2 opts_chunk$set(cache=TRUE,dpi=300,echo=FALSE,fig.width=7,  
3                 warning=FALSE,message=FALSE)  
4 @
```

- Great for journal articles, theses, books, ...

## 1 Reproducible Research

- Looks familiar ?
- Many Different Alternatives

## 2 R

- General Introduction
- Reproducible Documents: knitR
- **Introduction to R**
- Needful Packages by Hadley Wickam

# Data frames

A data frame is a data tables (with columns and rows). `mtcars` is a built-in data frame that we will use in the sequel

```
1 head(mtcars);
```

```
1           mpg  cyl  disp  hp  drat    wt    qsec  vs  am  gear  carb
2 Mazda RX4      21.0   6  160 110  3.90  2.620 16.46  0  1    4    4
3 Mazda RX4 Wag  21.0   6  160 110  3.90  2.875 17.02  0  1    4    4
4 Datsun 710     22.8   4  108  93  3.85  2.320 18.61  1  1    4    1
5 Hornet 4 Drive  21.4   6  258 110  3.08  3.215 19.44  1  0    3    1
6 Hornet Sportabout 18.7   8  360 175  3.15  3.440 17.02  0  0    3    2
7 Valiant        18.1   6  225 105  2.76  3.460 20.22  1  0    3    1
```

You can also load a data frame from a CSV file:

```
1 df <- read.csv("http://foo.org/mydata.csv", header=T,
2               strip.white=TRUE);
```

You will **get help** by using ?:

```
1 ?data.frame
2 ?rbind
3 ?cbind
```

# Exploring Content (1)

```
1 names(mtcars);
```

```
1 [1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am"  
2 [11] "carb"
```

```
1 str(mtcars);
```

```
1 'data.frame': 32 obs. of 11 variables:  
2 $ mpg : num 21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...  
3 $ cyl : num 6 6 4 6 8 6 8 4 4 6 ...  
4 $ disp: num 160 160 108 258 360 ...  
5 $ hp : num 110 110 93 110 175 105 245 62 95 123 ...  
6 $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...  
7 $ wt : num 2.62 2.88 2.32 3.21 3.44 ...  
8 $ qsec: num 16.5 17 18.6 19.4 17 ...  
9 $ vs : num 0 0 1 1 0 1 0 1 1 1 ...  
10 $ am : num 1 1 1 0 0 0 0 0 0 0 ...  
11 $ gear: num 4 4 4 3 3 3 3 4 4 4 ...  
12 $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

## Exploring Content (2)

```
1 dim(mtcars);  
2 length(mtcars);
```

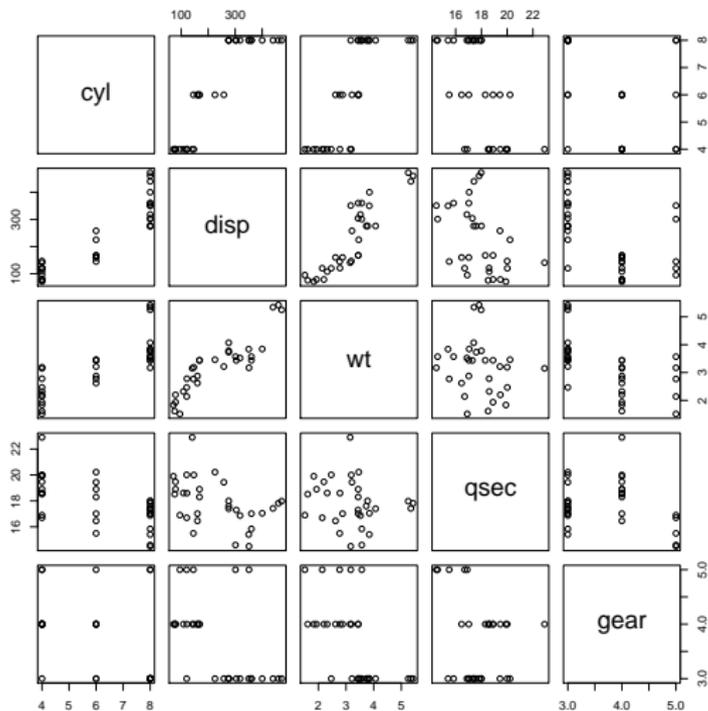
```
1 [1] 32 11  
2 [1] 11
```

```
1 summary(mtcars);
```

```
1           mpg           cyl           disp           hp  
2 Min.      :10.40   Min.      :4.000   Min.      : 71.1   Min.      : 52.0  
3 1st Qu.:15.43   1st Qu.:4.000   1st Qu.:120.8   1st Qu.: 96.5  
4 Median :19.20   Median :6.000   Median :196.3   Median :123.0  
5 Mean     :20.09   Mean     :6.188   Mean     :230.7   Mean     :146.7  
6 3rd Qu.:22.80   3rd Qu.:8.000   3rd Qu.:326.0   3rd Qu.:180.0  
7 Max.     :33.90   Max.     :8.000   Max.     :472.0   Max.     :335.0  
8           drat           wt           qsec           vs  
9 Min.      :2.760   Min.      :1.513   Min.      :14.50   Min.      :0.0000  
10 1st Qu.:3.080   1st Qu.:2.581   1st Qu.:16.89   1st Qu.:0.0000  
11 Median :3.695   Median :3.325   Median :17.71   Median :0.0000  
12 Mean     :3.597   Mean     :3.217   Mean     :17.85   Mean     :0.4375  
13 3rd Qu.:3.920   3rd Qu.:3.610   3rd Qu.:18.90   3rd Qu.:1.0000
```

# Exploring Content (3)

```
1 plot(mtcars[names(mtcars) %in% c("cyl","wt","disp","qsec","gear")])
```



# Accessing Content

```
1 mtcars$mpg
```

```
1 [1] 21.0 21.0 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 17.8 16.4 17
2 [16] 10.4 14.7 32.4 30.4 33.9 21.5 15.5 15.2 13.3 19.2 27.3 26.0 30
3 [31] 15.0 21.4
```

```
1 mtcars[2:5,]$mpg
```

```
1 [1] 21.0 22.8 21.4 18.7
```

```
1 mtcars[mtcars$mpg == 21.0,]
```

```
1           mpg cyl disp  hp drat   wt  qsec vs am gear carb
2 Mazda RX4    21   6  160 110  3.9 2.620 16.46 0  1   4    4
3 Mazda RX4 Wag 21   6  160 110  3.9 2.875 17.02 0  1   4    4
```

```
1 mtcars[mtcars$mpg == 21.0 & mtcars$wt > 2.7,]
```

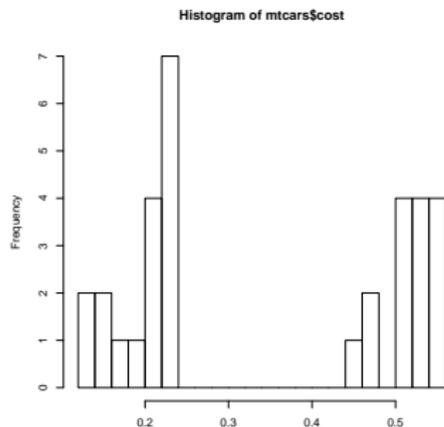
```
1           mpg cyl disp  hp drat   wt  qsec vs am gear carb
2 Mazda RX4 Wag 21   6  160 110  3.9 2.875 17.02 0  1   4    4
```

# Extending Content

```
1 mtcars$cost = log(mtcars$hp)*atan(mtcars$disp)/  
2   sqrt(mtcars$gear**5);  
3 mean(mtcars$cost);  
4 summary(mtcars$cost);
```

```
1 [1] 0.345994  
2   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
3 0.1261 0.2038 0.2353 0.3460 0.5202 0.5534
```

```
1 hist(mtcars$cost,breaks=20);
```



## 1 Reproducible Research

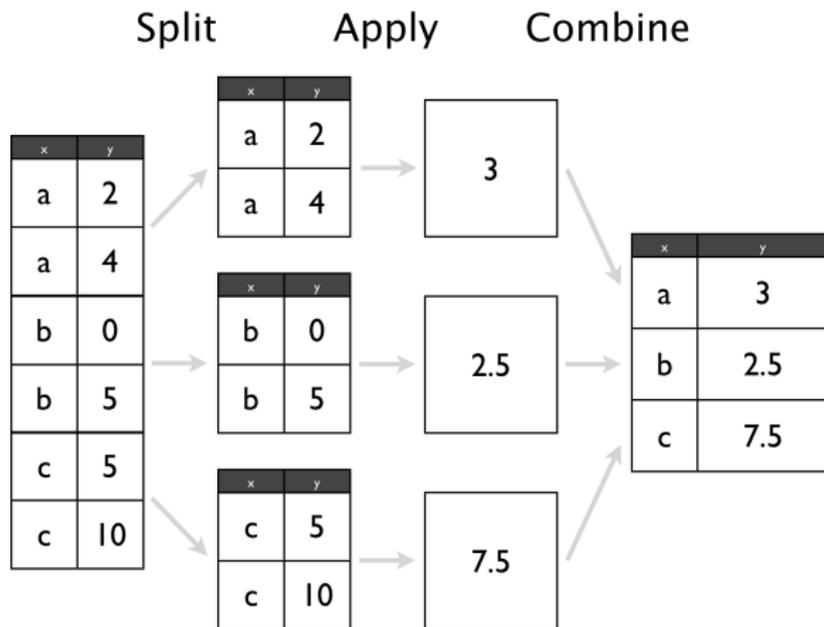
- Looks familiar ?
- Many Different Alternatives

## 2 R

- General Introduction
- Reproducible Documents: knitR
- Introduction to R
- **Needful Packages by Hadley Wickam**

# plyr: the Split-Apply-Combine Strategy

Have a look at <http://plyr.had.co.nz/09-user/> for a more detailed introduction.



## plyr: Powerfull One-liners

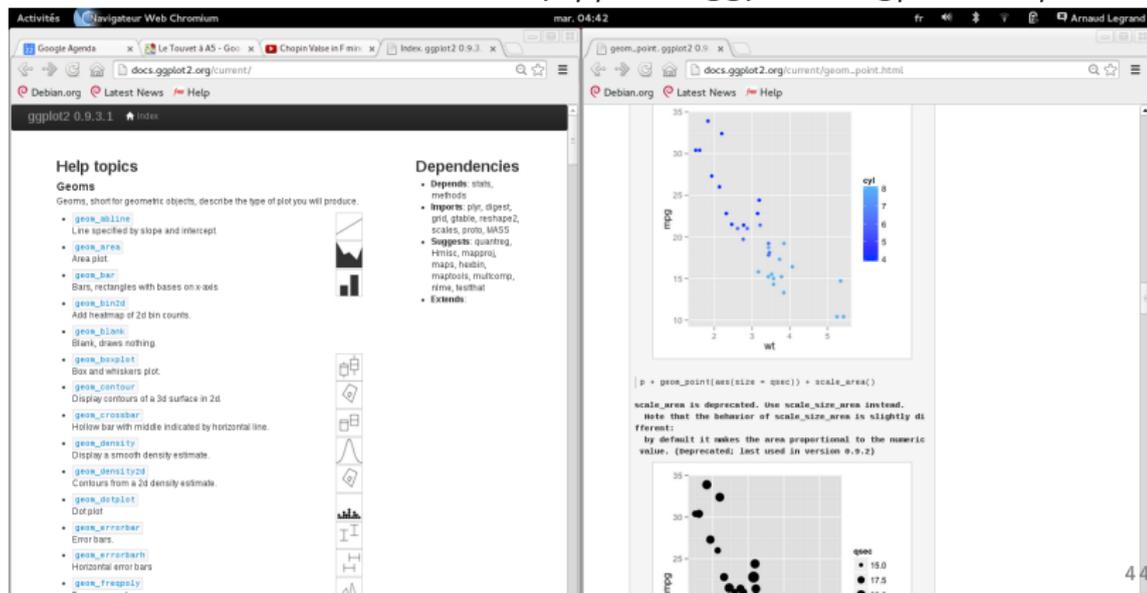
```
1 library(plyr)
2 mtcars_summarized = ddply(mtcars,c("cyl","carb"), summarize,
3     num = length(wt), wt_mean = mean(wt), wt_sd = sd(wt),
4     qsec_mean = mean(qsec), qsec_sd = sd(qsec));
5 mtcars_summarized
```

	cyl	carb	num	wt_mean	wt_sd	qsec_mean	qsec_sd
1	4	1	5	2.151000	0.2627118	19.37800	0.6121029
2	4	2	6	2.398000	0.7485412	18.93667	2.2924368
3	6	1	2	3.337500	0.1732412	19.83000	0.5515433
4	6	4	4	3.093750	0.4131460	17.67000	1.1249296
5	6	6	1	2.770000	NA	15.50000	NA
6	8	2	4	3.560000	0.1939502	17.06000	0.1783255
7	8	3	3	3.860000	0.1835756	17.66667	0.3055050
8	8	4	6	4.433167	1.0171431	16.49500	1.4424112
9	8	8	1	3.570000	NA	14.60000	NA

If your data is not in the right form **give a try to reshapeP/melt.**

# ggplot2: Modularity in Action

- ggplot2 builds on plyr and on a modular **grammar of graphics**
- **obnoxious** function with dozens of arguments
- **combine** small functions using layers and transformations
- **aesthetic** mapping between **observation characteristics** (data frame column names) and **graphical object variables**
- an incredible **documentation**: <http://docs.ggplot2.org/current/>



The screenshot shows the ggplot2 documentation website in a web browser. The page is titled "ggplot2 0.9.3.1" and has a navigation bar with "Index" and "Help". The main content is divided into two columns: "Help topics" and "Dependencies".

**Help topics**

**Geoms**  
Geoms, short for geometric objects, describe the type of plot you will produce.

- [geom\\_abline](#)  
Line specified by slope and intercept
- [geom\\_area](#)  
Area plot
- [geom\\_bar](#)  
Bars, rectangles with bases on x axis
- [geom\\_bin2d](#)  
Add heatmap of 2d bin counts.
- [geom\\_blank](#)  
Blank, draws nothing
- [geom\\_boxplot](#)  
Box and whiskers plot.
- [geom\\_contour](#)  
Display contours of a 3d surface in 2d
- [geom\\_crossbar](#)  
Hollow bar with middle indicated by horizontal line.
- [geom\\_density](#)  
Display a smooth density estimate.
- [geom\\_density2d](#)  
Contours from a 2d density estimate.
- [geom\\_dotplot](#)  
Dotplot
- [geom\\_errorbar](#)  
Error bars.
- [geom\\_errorbarh](#)  
Horizontal error bars
- [geom\\_freqpoly](#)

**Dependencies**

- **Depends:** stats, methods
- **Imports:** plyr, digest, grid, gtable, reshape2, scales, proto, MASS
- **Suggests:** quantreg, Hmisc, magrittr, maps, heatmap, mapproj, multcomp, rJava, testthat
- **Extends:**

**Code Snippets:**

```
p + geom_point(aes(size = qsec)) + scale_area()
```

**Text:**  
scale\_area is deprecated. Use scale\_size\_area instead.  
Note that the behavior of scale\_size\_area is slightly different:  
by default it makes the area proportional to the numeric value. (deprecated; last used in version 0.9.2)

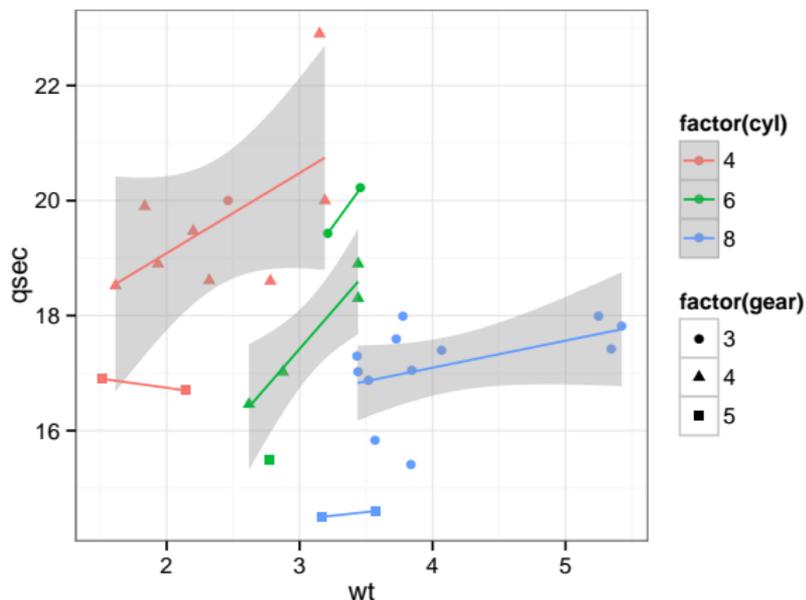
**Plots:**  
The top plot is a scatter plot of mpg vs wt, colored by cyl. The bottom plot is a scatter plot of mpg vs qsec, where the size of the points is proportional to qsec.





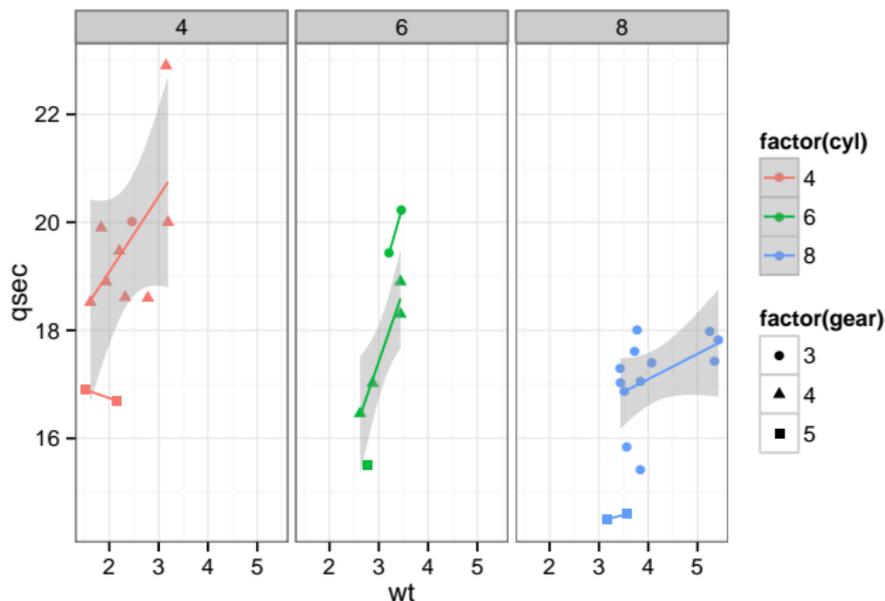
## ggplot2: Illustration (3)

```
1 ggplot(data = mtcars, aes(x=wt, y=qsec, color=factor(cyl),  
2   shape = factor(gear))) + geom_point() + theme_bw() +  
3   geom_smooth(method="lm");
```



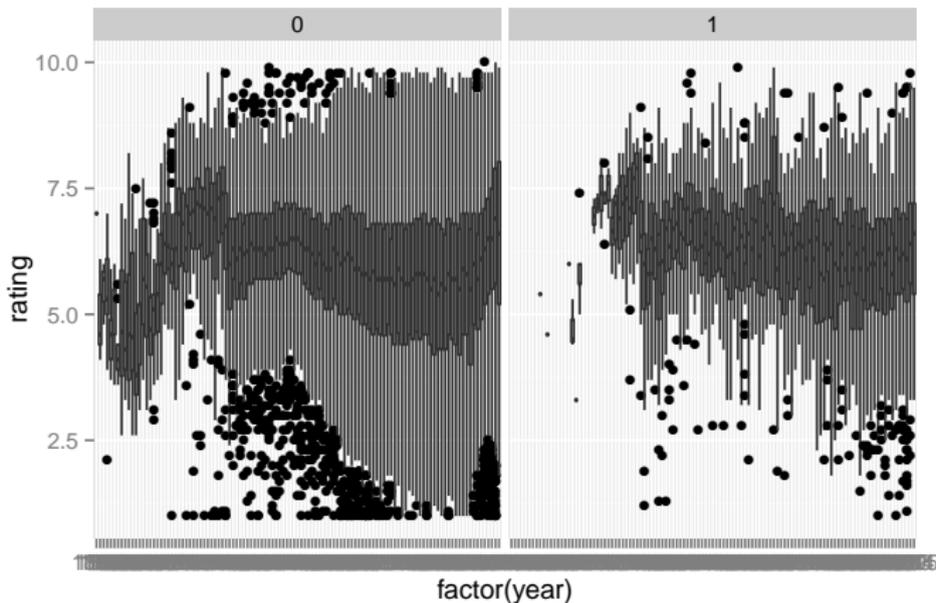
## ggplot2: Illustration (4)

```
1 ggplot(data = mtcars, aes(x=wt, y=qsec, color=factor(cyl),  
2   shape = factor(gear))) + geom_point() + theme_bw() +  
3   geom_smooth(method="lm") + facet_wrap(~ cyl);
```



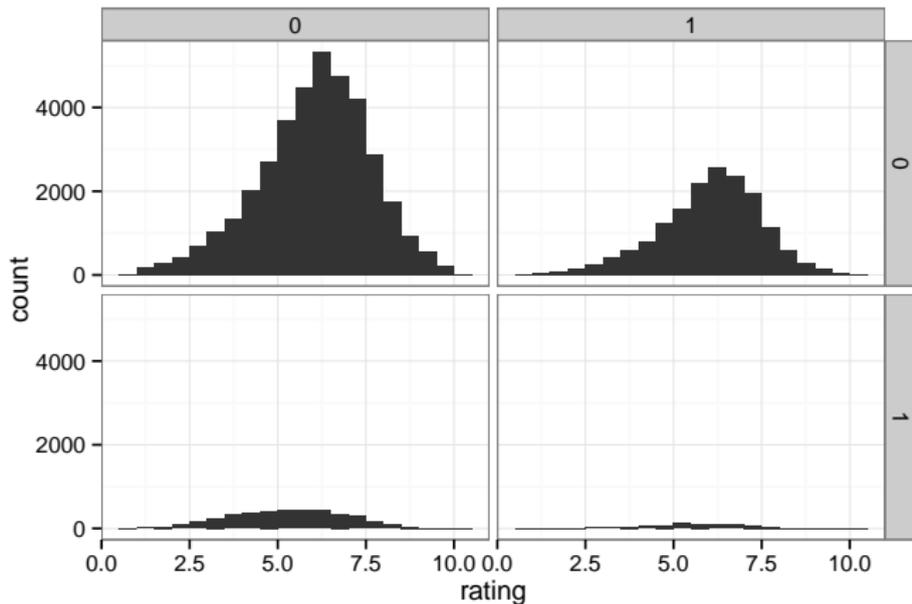
## ggplot2: Illustration (5)

```
1 ggplot(data = movies, aes(x=factor(year),y=rating)) +  
2   geom_boxplot() + facet_wrap(~Romance)
```



## ggplot2: Illustration (6)

```
1 ggplot(movies, aes(x = rating)) + geom_histogram(binwidth = 0.5)+  
2   facet_grid(Action ~ Comedy) + theme_bw();
```



# Take away Message

- R is a great tool but is only a tool. There is no magic. You need to understand what you are doing and get a minimal training in statistics.
- It is one of the building block of **reproducible research** (the *reproducible analysis* block) and **will save you a lot of time**.
- Read at least Jain's book: *The Art of Computer Systems Performance Analysis*.
- Jean-Marc Vincent and myself give a **set of tutorials on performance evaluation** in M2R:

[http://mescal.imag.fr/membres/arnaud.legrand/teaching/2013/M2R\\_EP.php](http://mescal.imag.fr/membres/arnaud.legrand/teaching/2013/M2R_EP.php)

- There are interesting **online courses** on coursera
  - <https://www.coursera.org/course/compdata>
  - <https://www.coursera.org/course/repdata>

## About these slides

They have been composed in `org-mode` and generated with `emacs`, `beamer`, and `pyglist/pygments` for the pretty printing.