

# R crash course

Arnaud Legrand and Jean-Marc Vincent

Scientific Methodology and Performance Evaluation  
ENS Lyon, November 2016

# Outline

## ① R/knitr Crash Course

- General Introduction

- Reproducible Documents: knitr

- Introduction to R

## ② Needful R Packages by Hadley Wickam

- Plyr And Dplyr

- Ggplot2

- Reshape and tidyR

- Now let's play!

- Conclusion

# 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 such 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 3.2.0 (2015-04-16) -- "Full of Ingredients"
2 Copyright (C) 2015 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.
```

# Install a few cool packages

R has it's own package management mechanism so just run R and type the following commands:

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

```
1 install.packages("plyr")
2   # or better: install.packages("dplyr")
3 install.packages("reshape")
4   # or better; install.packages("tidyr")
5 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, I will briefly show you **rstudio** (<https://www.rstudio.com/>) and later how to use org-mode



# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

Reshape and tidyR

Now let's play!

Conclusion

# Rstudio screenshot

RStudio

File Edit Code View Project Workspace Plots Tools Help

Go to file/function

Project: (None)

Workspace History

Load Save Import Dataset Clear All

Data

df 10 obs. of 2 variables

Values

x integer[10]

y numeric[10]

```
28
29- ##{r basicconsole}
30 x <- 1:10
31 y <- round(rnorm(10, x, 1), 2)
32 df <- data.frame(x, y)
33 df
34
35
36 ## Plots
37 Images generated by 'knitr' are saved in a figures folder. However,
38 | they also appear to be represented in the HTML output using a [data
39 | URI scheme]( http://en.wikipedia.org/wiki/Data_URI_scheme). This
40 | means that you can paste the HTML into a blog post or discussion
41 | forum and you don't have to worry about finding a place to store the
42 | images; they're embedded in the HTML.
43
44 ### Simple plot
45 Here is a basic plot using base graphics:
46
47 ##{r simpleplot}
48 plot(x)
49
50 ##{r simpleplot}
51 plot(x)
```

Console

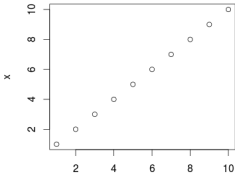
~/research/statistics/markdown-meetup-2012/

'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

```
> set.seed(1234)
> library(ggplot2)
> library(lattice)
> x <- 1:10
> y <- round(rnorm(10, x, 1), 2)
> df <- data.frame(x, y)
> 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
> plot(x)
```

Files Plots Packages Help

Zoom Export Clear All



Index

# 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, ...

# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

Reshape and tidyR

Now let's play!

Conclusion

# 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
```

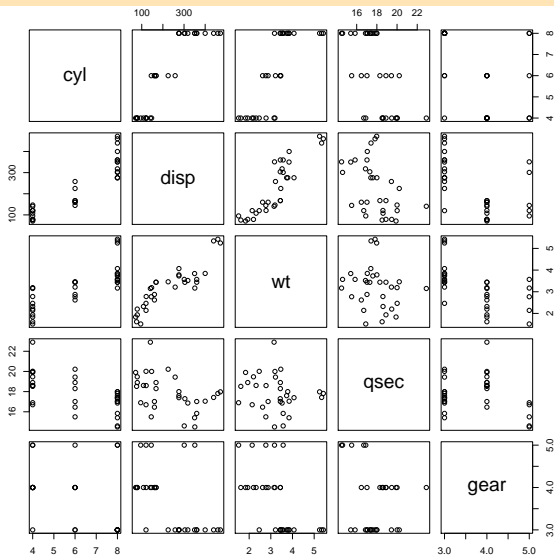
```
1 summary(mtcars);
```

	mpg	cyl	disp	hp
1				
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 17.8
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 30.4
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
```

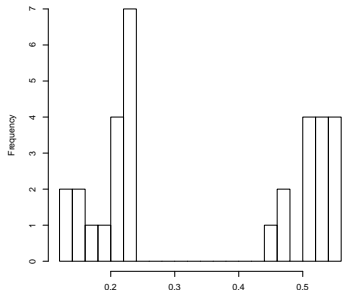
# Extending Content

```
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);
```

Histogram of mtcars\$cost



# 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**
- It provides you an access to any statistical method you ever dreamt of
- Read at least Jain's book: **The Art of Computer Systems Performance Analysis**
- There are introductory **online courses** (from John Hopkins university) on coursera which you may want to follow

# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

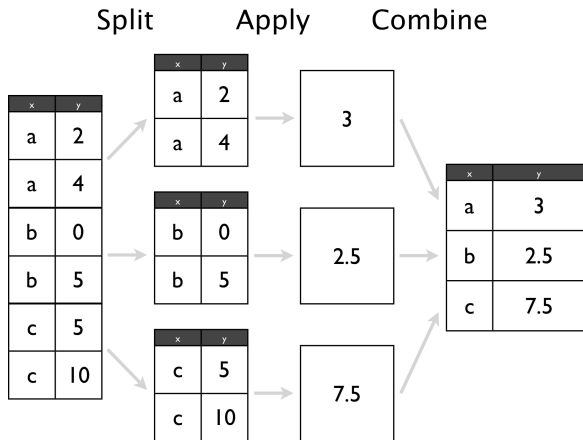
Reshape and tidyR

Now let's play!

Conclusion

## plyr: the Split-Apply-Combine Strategy

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



# plyr: Powerful 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	1	4	1	5	2.151000	0.2627118	19.37800	0.6121029
2	2	4	2	6	2.398000	0.7485412	18.93667	2.2924368
3	3	6	1	2	3.337500	0.1732412	19.83000	0.5515433
4	4	6	4	4	3.093750	0.4131460	17.67000	1.1249296
5	5	6	6	1	2.770000	NA	15.50000	NA
6	6	8	2	4	3.560000	0.1939502	17.06000	0.1783255
7	7	8	3	3	3.860000	0.1835756	17.66667	0.3055050
8	8	8	4	6	4.433167	1.0171431	16.49500	1.4424112
9	9	8	8	1	3.570000	NA	14.60000	NA

## plyr next generation = dplyr

It's much much faster and more readable. The *tutorial* is great...

```

1 library(dplyr)
2 mtcars %>% group_by(cyl,carb) %>%
3     select(wt,qsec) %>%
4     summarise(num = n(),
5     wt_mean = mean(wt), wt_sd = sd(wt),
6     qsec_mean = mean(qsec), qsec_sd = sd(qsec)) %>%
7     filter(num>=1)

```

```

1 Source: local data frame [9 x 7]
2 Groups: cyl
3
4   cyl carb num  wt_mean    wt_sd qsec_mean  qsec_sd
5 1    4    1   5 2.151000 0.2627118  19.37800 0.6121029
6 2    4    2   6 2.398000 0.7485412  18.93667 2.2924368
7 3    6    1   2 3.337500 0.1732412  19.83000 0.5515433

```



# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

Reshape and tidyR

Now let's play!

Conclusion

# 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/>

Activities | Navigateur Web Chromium | mar. 04:42 | fr | Arnaud Legrand

Google Agenda | Le Touvet à A5 - Geo | Chopin Valse in F min | Index: ggplot2 0.9.3 | docs.ggplot2.org/current/ | Debian.org | Latest News | Help

## ggplot2 0.9.3.1

Index

### 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

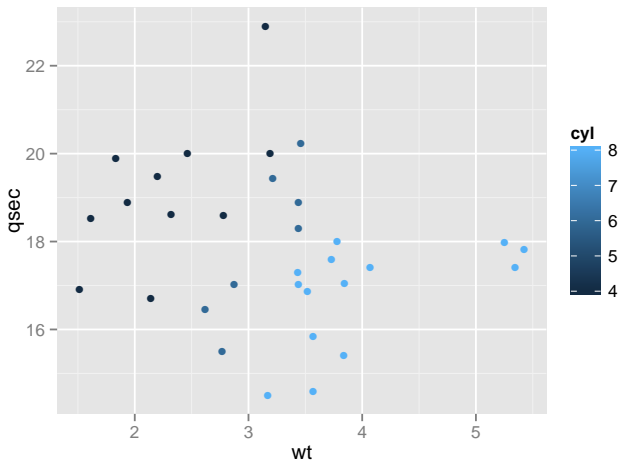
- **Depends:** stats, methods
- **Imports:** plyr, digest, grid, reshape2, scales, proto, MASS
- **Suggests:** quantreg, Rmisc, maptools, maps, heatmap, maptools, multcomp, rma, lme4
- **Extends:**

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

`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)

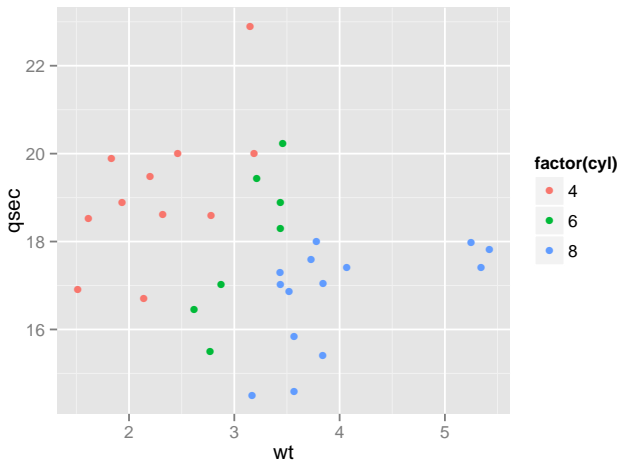
## ggplot2: Illustration (1)

```
1 ggplot(data = mtcars, aes(x=wt, y=qsec, color=cyl)) +  
2   geom_point();
```



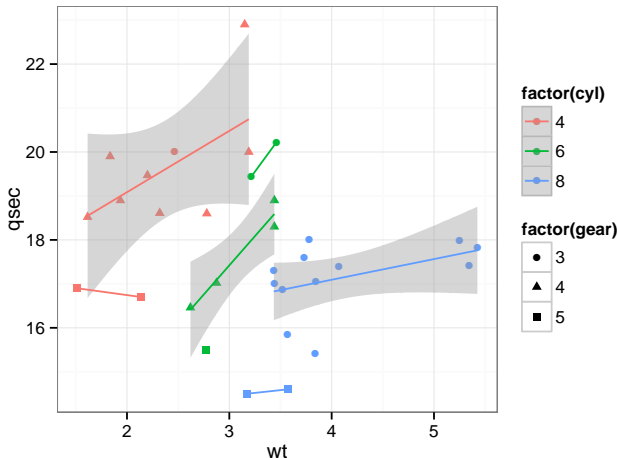
## ggplot2: Illustration (2)

```
1 ggplot(data = mtcars, aes(x=wt, y=qsec, color=factor(cyl))) +  
2   geom_point();
```



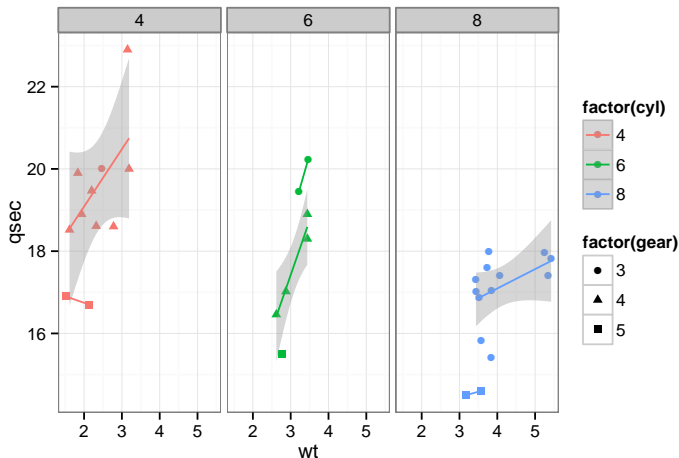
## ggplot2: Illustration (3)

```
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");
```



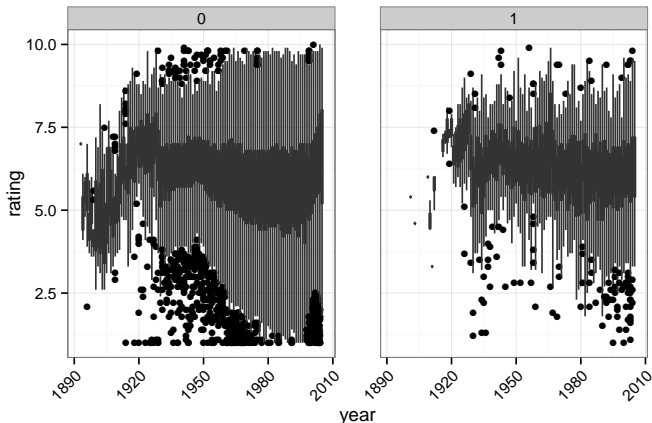
## ggplot2: Illustration (4)

```
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);
```



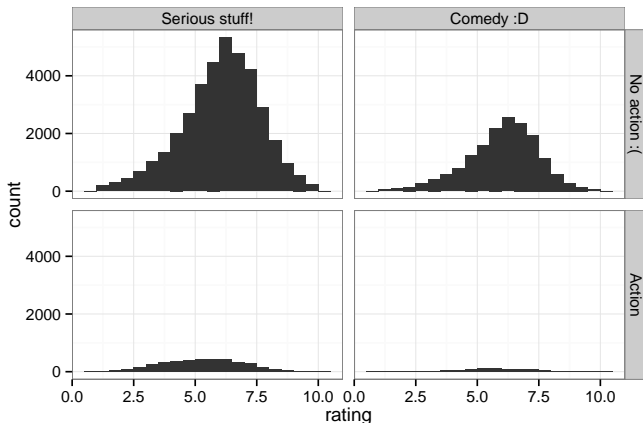
## ggplot2: Illustration (5)

```
1 ggplot(data = movies, aes(x=year,y=rating,group=factor(year))) +  
2   geom_boxplot() + facet_wrap(~Romance) + theme_bw() +  
3   theme(axis.text.x = element_text(angle = 45, hjust = 1),  
4         panel.margin = unit(2, "lines"));
```



## ggplot2: Illustration (6)

```
1 ggplot(movies, aes(x = rating)) + geom_histogram(binwidth = 0.5) +  
2   facet_grid(Action ~ Comedy, labeller=mf_labeller) +  
3   theme_bw() + theme(panel.margin = unit(.5, "lines"));
```





# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

Reshape and tidyR

Now let's play!

Conclusion

# "Messy" data

When using ggplot or plyr, your data may not in the right shape, in which case you should **give a try to reshape/melt**

```
1 messy <- data.frame(  
2   name = c("Wilbur", "Petunia", "Gregory"),  
3   a = c(67, 80, 64),  
4   b = c(56, 90, 50)  
5 )  
6 messy
```

```
1      name  a  b  
2 1 Wilbur 67 56  
3 2 Petunia 80 90  
4 3 Gregory 64 50
```

- a and b are two different types of drugs and the values correspond to heart rate
- ggplot faceting or coloring based on the drug type is a pain
- we need a way to make "wide" data longer

# Reshape

```
1 library(reshape)
2 cleaner = melt(messy, c("name"))
3 names(cleaner)=c("name", "drug", "heartrate")
4 cleaner
```

```
1      name drug heartrate
2 1 Wilbur    a         67
3 2 Petunia   a         80
4 3 Gregory   a         64
5 4 Wilbur    b         56
6 5 Petunia   b         90
7 6 Gregory   b         50
```

# Tidyr

Just like `plyr`, `reshape` is a little magical. `tidyr` is the new generation (faster, more coherent). Again, the *tutorial* is great.

```
1 library(tidyr)
2 library(dplyr)
3 messy %>% gather(drug, heartrate, -name)
```

```
1      name drug heartrate
2 1 Wilbur   a      67
3 2 Petunia  a      80
4 3 Gregory  a      64
5 4 Wilbur   b      56
6 5 Petunia  b      90
7 6 Gregory  b      50
```

**Hint:** Avoid mixing old-generation with new-generation as it overrides some function names and leads to weird behaviors

# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

Reshape and tidyR

Now let's play!

Conclusion

# Summarizing information

You may like these **cheat sheets**:

<https://www.rstudio.com/resources/cheatsheets/>

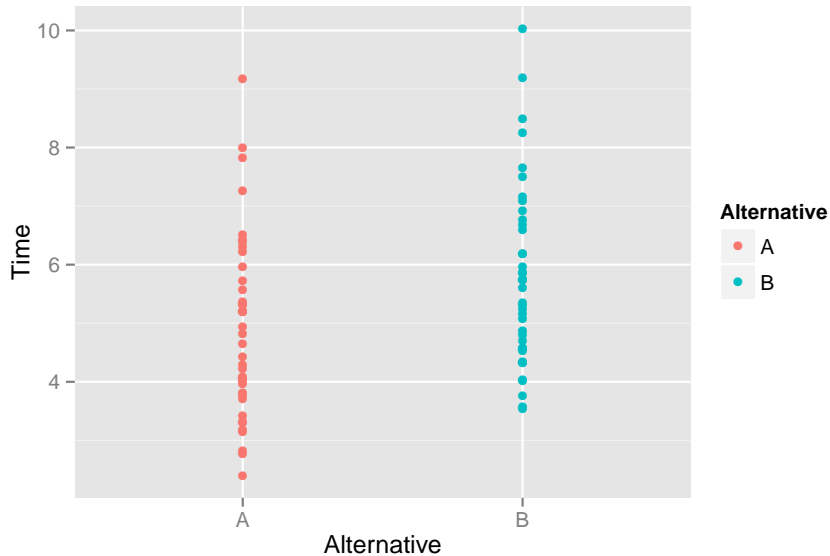
```
1 df = read.csv("data/set1.csv",header=T)
2 # Alternatively: read.csv("https://raw.githubusercontent.com/
3 #                       aleggrand/SMPE/master/lectures/data/set1.csv")
4 head(df,n=2)
```

```
1           A           B
2 1 7.256717 8.261171
3 2 3.813100 4.335301
```

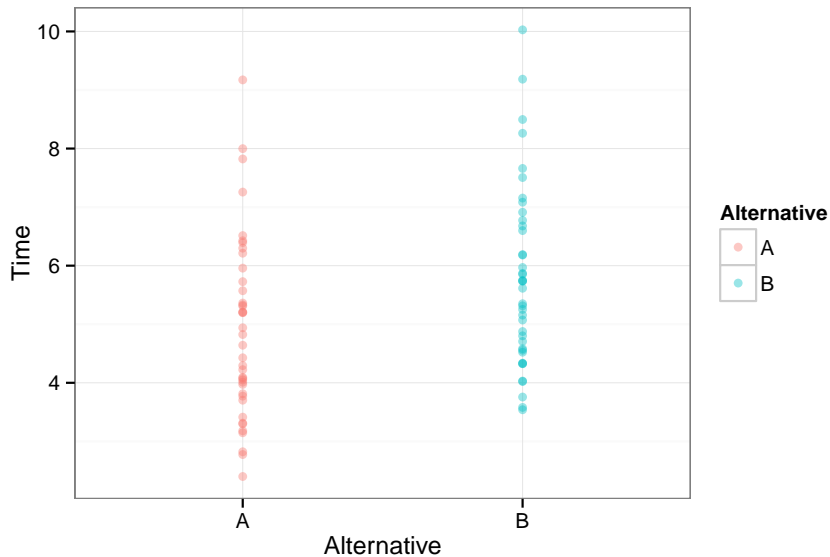
Get the following summary using `plyr/reshape` or `dplyr/tydir`:

```
1 Source: local data frame [2 x 6]
2
3   Alternative num    mean      sd    min    max
4 1           A  40 4.903817 1.544423 2.400016 9.172525
5 2           B  40 5.783643 1.542987 3.539874 10.027147
```

# Plot the data

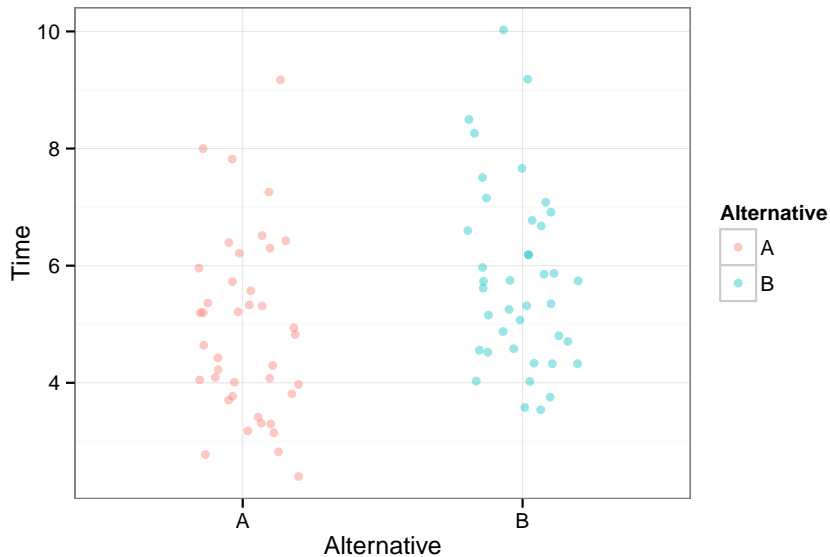


# Alleviate over-plotting

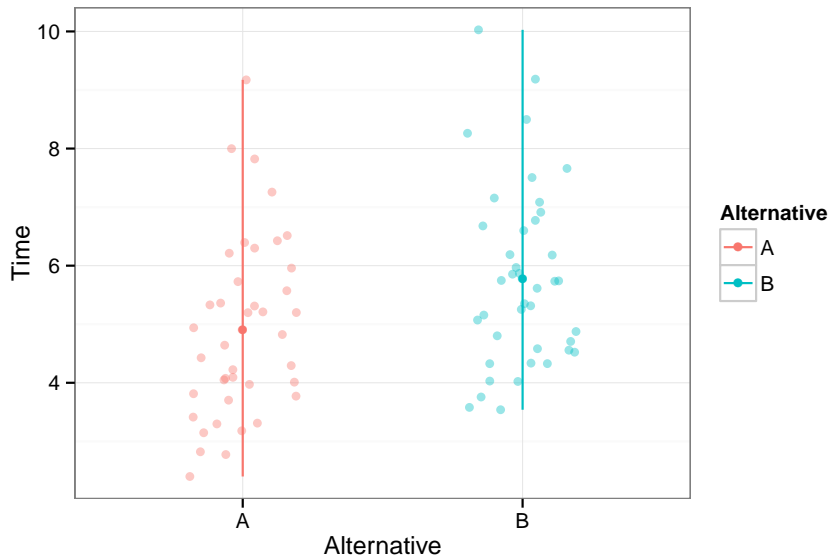




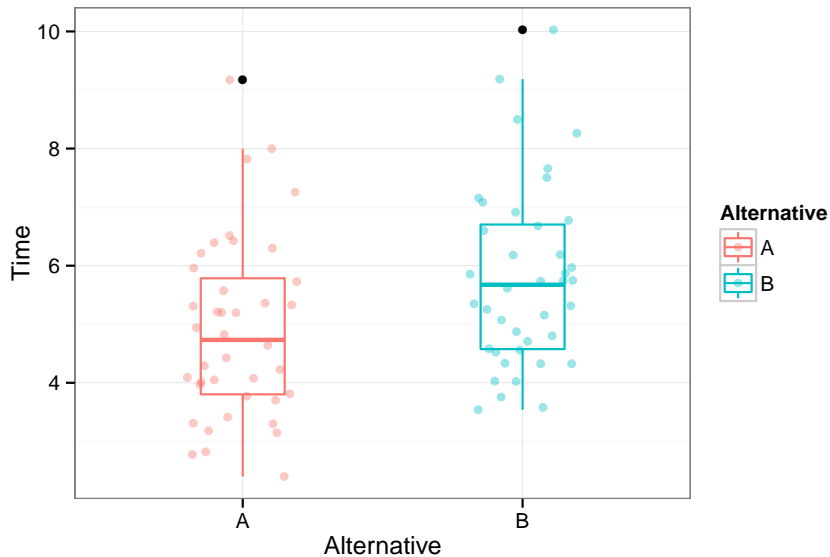
# Avoid over-plotting



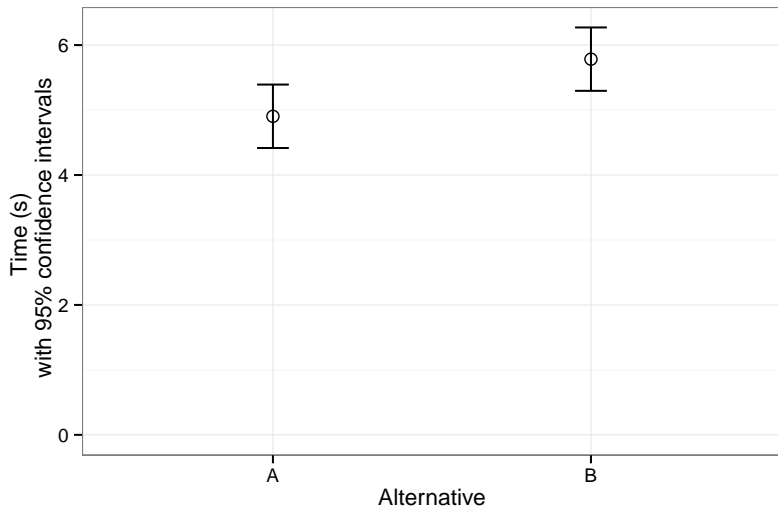
## Add summary information



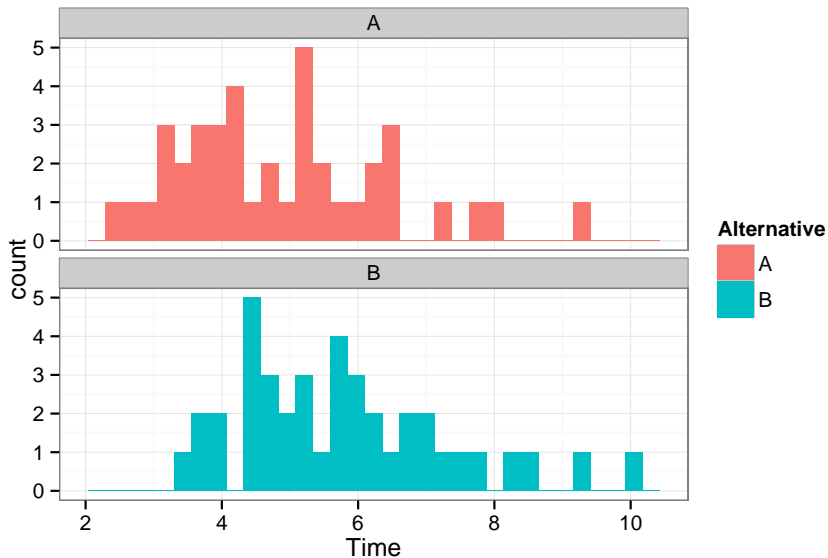
## Add more standard summaries



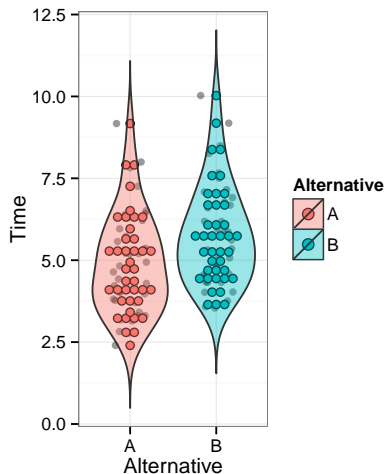
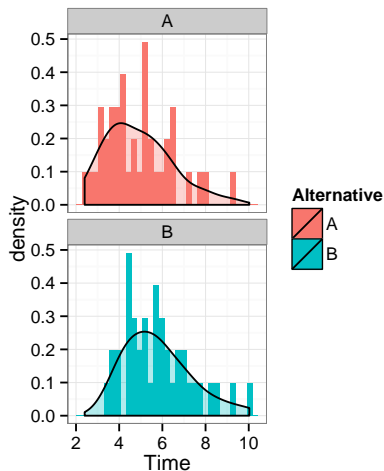
## Or depict confidence intervals



Or use histograms. . .



# Be careful with fancy plots you do not fully understand!



# Outline

## ① R/knitr Crash Course

General Introduction

Reproducible Documents: knitr

Introduction to R

## ② Needful R Packages by Hadley Wickam

Plyr And Dplyr

Ggplot2

Reshape and tidyR

Now let's play!

Conclusion

# Take away Message

- R, ggplot and other such tools are **incredibly powerful for presenting data**. They are much more high level than any other tools I have seen so far.
- Mastering it **will save you a lot of time** as it will allow to look at your data through **different angles** and thus **check many hypothesis** and **present them in the best possible way**
- Read at least Jain's book: **The Art of Computer Systems Performance Analysis**
- However, you may have started understanding that a visualization is meant to check or to illustrate one particular aspect and that this "aspect" relies on a **mathematical model**. I will thus explain you in the next lecture what this model is.

**To do for the Next Time:** Use what you just learned to improve your data analysis, the article you're currently writing, ...