



# STT 1000 - STATISTIQUES

ARTHUR CHARPENTIER



# “Logiciels de Statistique”

cf [comparaison des logiciels de statistiques](#) sur wikipedia

La connaissance d'un logiciel de programmation/calculs/analyse de données est devenu un élément incontournable pour les chercheurs, praticiens, industriels et donc des étudiants.



Pourquoi R ? Notamment parce qu'il est gratuit, libre, multiplateforme et bénéficie de la plus grande communauté au monde.

# Langage S/Logiciel R/RStudio

**S** “*est un langage de programmation de très haut niveau et un environnement d'analyse des données et des graphiques conçu dans les années 1975-1976 par John Chambers. Les deux interpréteurs modernes de S sont S-Plus et R*”

**R** “*naît en 1993 comme un projet de recherche de Ross Ihaka et Robert Gentleman à l'université d'Auckland, sous une licence GNU GPL*” (logiciel libre et gratuit)

**RStudio**: Logiciel libre, multi-plateforme, gratuit ; Interface graphique permettant un environnement de développement intégré (IDE).

Un IDE n'est pas une interface graphique au sens de SPSS ou SAS, qui permettrait d'utiliser le logiciel à travers des menus et des boîtes de dialogue : il s'agit d'un environnement facilitant la saisie, l'exécution de code, la visualisation des résultats, etc.

# Foire aux Questions

- ▶ *Ai-je besoin d'installer les logiciels R et RStudio sur une machine personnelle ?*

**Oui**, ça sera plus simple / **Non**, il existe de versions en ligne (<https://rstudio.cloud/>)

- ▶ *Ai-je besoin de maîtriser parfaitement le logiciel R pour réussir le cours avec succès ?*

**Non**, la maîtrise de R n'est pas une fin en soi (dans ce cours). R est juste un outils (très pratique)

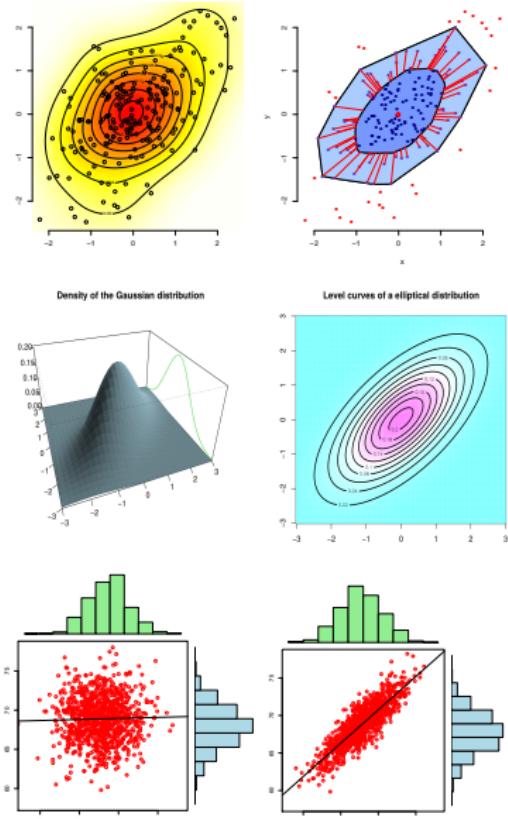
- ▶ *R/RStudio sera-t-il utilisé comme outil d'évaluation ?*

**Oui**, les examens pourront contenir des sorties logicielles qu'il conviendra d'interpréter,

# S/R Language

## Statistiques et graphiques

```
1 > x = c(0,1,2,3,4,5,6,7,8,9,10)
2 > mean(x)
3 [1] 5
4 > quantile(x, .95)
5 95%
6 9.5
7 > quantile(x, .95, type=1)
8 95%
9 10
10 > quantile(x, .95, type=4)
11 95%
12 9.45
13 > quantile(x, .95, type=5)
14 95%
15 9.95
```



# RStudio

The screenshot shows the RStudio interface with the following components:

- Top Bar:** Contains tabs for "Untitled1", "Source on Save", "Run", "Source", and "Project: (None)".
- Environment Tab:** Shows the global environment with various objects listed:
  - A: 10 obs. of 2 variables
  - acp: List of 5
  - anova: List of 13
  - AOV: List of 13
  - apartments: 1000 obs. of 6 variables  
ALTREIC: num [1:45, 1:3] -62.7 807 907.8 -53.7 213...
  - B: 7 obs. of 2 variables
  - b2: 28 obs. of 5 variables
  - base: 891 obs. of 8 variables
  - BT: List of 2
  - C: num [1:3, 1:558] 0.1677 -0.0525 0.3725 0...
  - coh: List of 7
  - cors: num [1:4, 1:4] 1 0.503 0.856 0.493 0.503 -
  - d: 28 obs. of 2 variables
  - dat14boot: List of 11
  - database: 6 obs. of 32 variables
  - Davis: 200 obs. of 5 variables
  - dc: 12 obs. of 12 variables
  - dd: 50 obs. of 5 variables
- Console Tab:** Displays the R startup message and workspace details.

```
R version 3.6.3 (2020-02-29) -- "Holding the Windsock"
Copyright (C) 2020 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

  Natural language support is running in an English locale

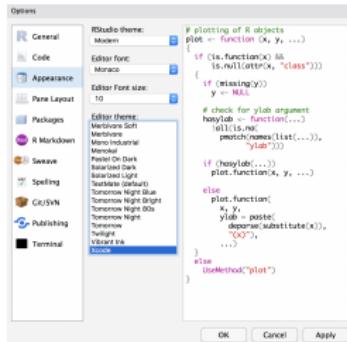
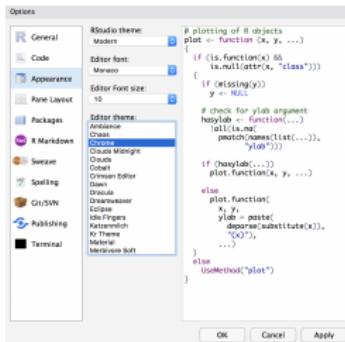
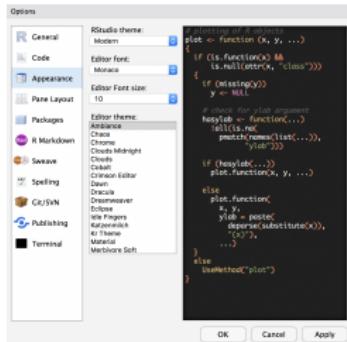
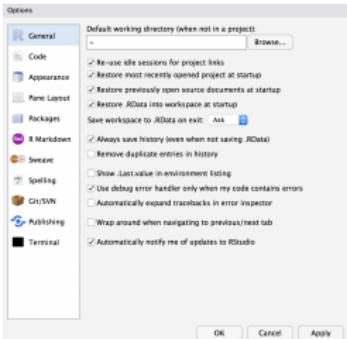
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/RData]
```
- Bottom Status Bar:** Shows navigation icons for files, plots, packages, help, viewer, and export.

# RStudio: appearance

Ma version de RStudio est sombre  
on peut la paramétrer...

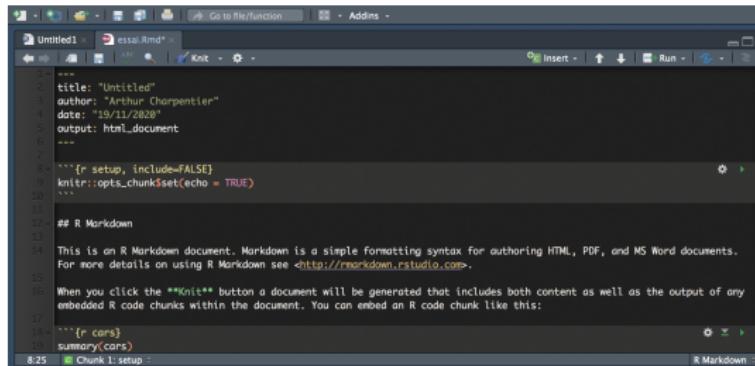


# RStudio

The screenshot shows the RStudio interface with the following components:

- Code Editor:** Displays an R Markdown document titled "Untitled1" with code for setting chunk options and a summary section.
- Environment Browser:** Shows the global environment with various objects listed, such as A, acp, anova, ADV, apartments, AUTREIC, B, b2, base, BT, c, coh, cors, d, dat4boot, database, Davis, dc, and dd.
- Help Viewer:** Displays the "t.test" function documentation under "Student's t-Test".

# RStudio



The screenshot shows the RStudio interface with the code editor window open. The file is named "essai.Rmd". The code in the editor is:

```
1 title: "Untitled"
2 author: "Arthur Charpentier"
3 date: "19/12/2020"
4 output: html_document
5
6
7 ````{r setup, include=FALSE}
8 knitr::opts_chunk$set(echo = TRUE)
9 ```
10
11 ## R Markdown
12
13 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents.
14 For more details on using R Markdown see <http://rmarkdown.rstudio.com>.
15
16 When you click the **Knit** button a document will be generated that includes both content as well as the output of any
17 embedded R code chunks within the document. You can embed an R code chunk like this:
18
19 ````{r cars}
20 summary(cars)
21 ````
```

The status bar at the bottom indicates "8:25" and "Chunk 1: setup".

La fenêtre **editor** (haut à gauche)

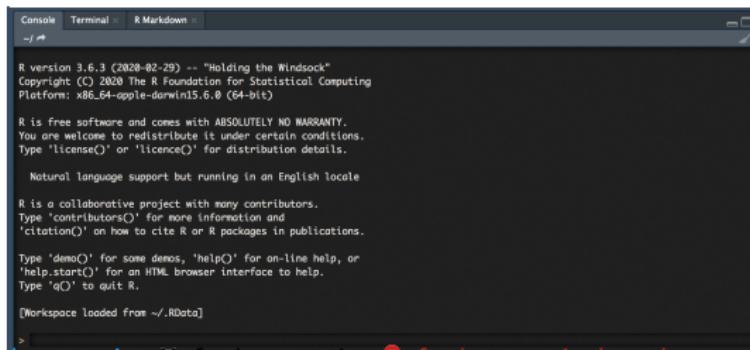
On peut y taper une succession de commandes (que l'on pourra exécuter ensuite)

```
1 # 5 observations
2 v = c(0.89367, -1.04729, 1.97133, -0.38363, 1.65414)
3 mean(v)
4 f = function(x) sum((v-x)^2)
5 optim(0, f)
```

# RStudio

La fenêtre **console** (bas à gauche) : après avoir tapé 'run',

```
1 > v = c(0.89367, -1.04729, 1.97133, -0.38363, 1.65414)
2 > mean(v)
3 [1] 0.617644
4 > f = function(x) sum((v-x)^2)
5 > optim(0, f)
6 $par
7 [1] 0.6175781
8 Warning message:
9 In optim(0, f) : one-dimensional optimization by
   Nelder-Mead is unreliable:
10 use "Brent" or optimize() directly
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The window displays the standard R startup messages, including the version number (R version 3.6.3), copyright information, and instructions for redistribution. It also mentions natural language support and collaborative project details. At the bottom, it indicates that the workspace was loaded from a specific directory.

```
Console Terminal R Markdown
~/R

R version 3.6.3 (2020-02-29) -- "Holding the Windsock"
Copyright (C) 2020 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/RData]
```

La fenêtre **environment** (haut à droite)  
donne la liste des objets en mémoire

```
1 v      num [1:5] 0.894 -1.047...
2 f      function(x)
```

The screenshot shows the RStudio interface with the 'Environment' tab selected in the top navigation bar. The main pane displays a list of objects currently stored in memory, categorized under 'Data'. Each object is listed with its name, type, and a brief description of its contents. For example, 'A' is a '10 obs. of 2 variables' data frame containing numerical values. Other objects listed include 'acp', 'anova', 'AGV', 'apartments', 'AUTREIC', 'B', 'b2', 'base', 'BT', 'C', 'cah', 'cars', 'd', 'dot14boot', 'database', 'Davis', 'dc', and 'rid'. Each entry has a small circular icon to its left and a 'Search' icon to its right.

Object	Type / Description
A	10 obs. of 2 variables
acp	List of 5
anova	List of 13
AGV	List of 13
apartments	1000 obs. of 6 variables
AUTREIC	num [1:45, 1:3] -62.7 807 907.8 -53.7 213...
B	7 obs. of 2 variables
b2	20 obs. of 5 variables
base	891 obs. of 8 variables
BT	List of 2
C	num [1:3, 1:558] 0.1677 -0.0525 0.3725 0...
cah	List of 7
cars	num [1:4, 1:4] 1 0.503 0.856 0.493 0.503 ...
d	20 obs. of 2 variables
dot14boot	List of 11
database	6 obs. of 32 variables
Davis	200 obs. of 5 variables
dc	12 obs. of 12 variables
rid	50 obs. of 5 variables

# RStudio

## La fenêtre packages (bas à droite)

```
1 > install.packages("vcd")
```



[CRAN](#)  
[Mirrors](#)  
[What's new?](#)  
[Task Views](#)  
[Search](#)

[About R](#)  
[R Homepage](#)  
[The R Journal](#)

[Software](#)  
[R Sources](#)  
[R Binaries](#)  
[Packages](#)  
[Other](#)

[Documentation](#)  
[Manuals](#)  
[FAQs](#)  
[Contributed](#)

### Contributed Packages

#### Available Packages

Currently, the CRAN package repository features 17628 available packages.

[Table of available packages, sorted by date of publication](#)

[Table of available packages, sorted by name](#)

#### Installation of Packages

Please type `help("INSTALL")` or `help("install.packages")` in R for information on how to install packages from this repository. The manual [R Installation and Administration](#) (also contained in the R base sources) explains the process in detail.

[CRAN Task Views](#) allow you to browse packages by topic and provide tools to automatically install all packages for special areas of interest. Currently, 41 views are available.

#### Package Check Results

All packages are tested regularly on machines running [Debian GNU/Linux](#), [Fedora](#), macOS (formerly OS X), Solaris and Windows.

The results are summarized in the [check summary](#) (some [timings](#) are also available). Additional details for Windows checking and building can be found in the [Windows check summary](#).

Name	Description	Version
<b>System Library</b>		
abind	Combine Multidimensional Arrays	1.4-5
acepack	ACE and AVAS for Selecting Multiple Regression Transformations	1.4.1
AER	Applied Econometrics with R	1.2-9
ape	Analyses of Phylogenetics and Evolution	5.4
arm	Data Analysis Using Regression and Multilevel/Hierarchical Models	1.11-1
ash	David Scott's ASH Routines	1.0-15
askpass	Safe Password Entry for R, Git, and SSH	1.1
assertthat	Easy Pre and Post Assertions	0.2.1
backports	Reimplementations of Functions Introduced Since R 3-0.0	1.1.7
base64enc	Tools for base64 encoding	0.1-3
BenfordTests	Statistical Tests for Evaluating Conformity to Benford's Law	1.2.0
bit	Bit-Oriented Memory Efficient Tools	1.22.0

## Sample Quantiles

### Description

The generic function `quantile` produces sample quantiles corresponding to the given probabilities. The smallest observation corresponds to a probability of 0 and the largest to a probability of 1

### Usage

```
1 > quantile(x, probs = seq(0, 1, 0.25), na.rm = FALSE,  
  names = TRUE, type = 7, ...)
```

### Arguments

`x`: numeric vector whose sample quantiles are wanted, or an object of a class for which a method has been defined (see also 'details'). NA and NaN values are not allowed in numeric vectors unless `na.rm` is TRUE

### References

Hyndman, R. J. and Fan, Y. (1996) Sample quantiles in statistical packages, *American Statistician* 50, 361–365. doi: 10.2307/2684934.

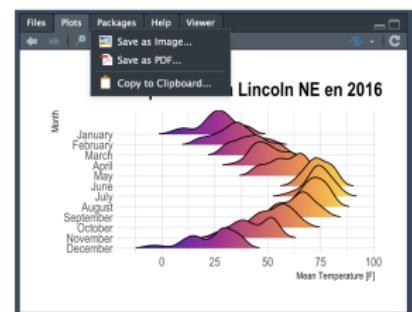
# RStudio: graphiques

Pour les graphiques (de base)

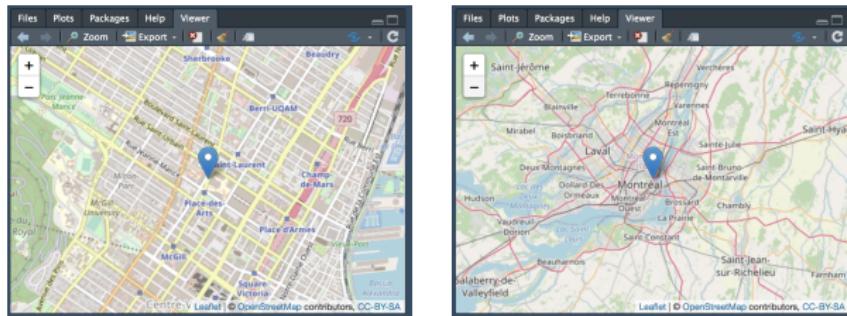
1. il faut créer la fenêtre graphique (avec la fonction plot)

```
1 > plot(cars)
```

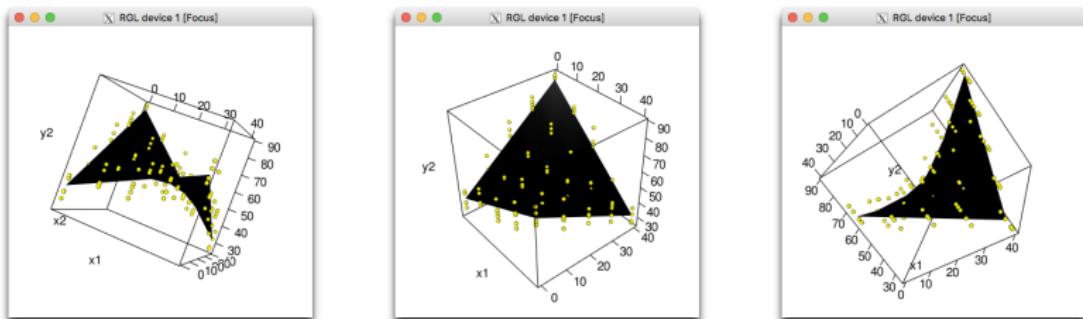
2. on peut alors superposer des courbes, des points, etc



Certains graphiques (dynamiques) apparaîtront dans la fenêtre viewer (e.g. leaflet)



ou dans une fenêtre extérieure (`rgl`)



# Fichier STT1000.RData

Les fichiers **.Rdata** contiennent des données (format lisible en R).

```
1 > url = "http://freakonometrics.free.fr/STT1000.RData"
2 > download.file(url,"STT1000.RData")
3 trying URL
4 Content type 'text/plain' length 459595 bytes (448 KB)
5 =====
6 downloaded 448 KB
7 > load("STT1000.RData")
```



# Fichier STT1000.RData

```
1 > ls()
2 [1] "a"                      "alcool"                  "aspirine"
3 [4] "attente"                "b"                      "babies"
```

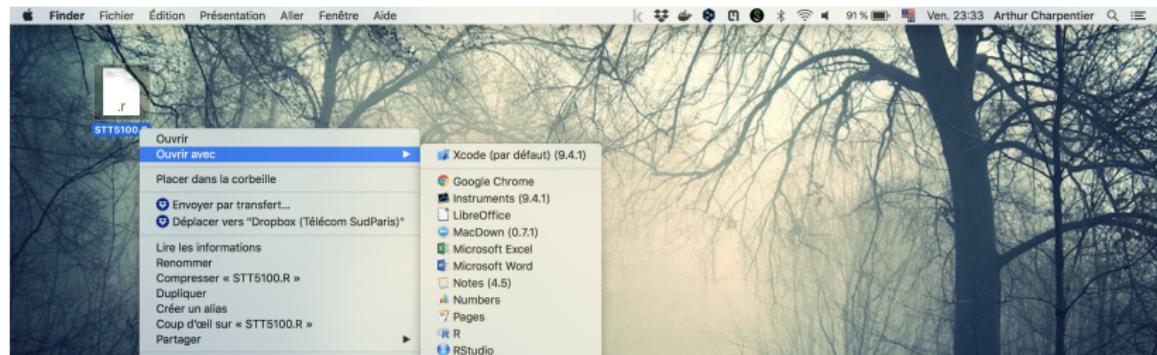
(etc)

```
1 > tail(Davis)
2   sex weight height reportedWeight reportedHeight
3 195   F     62    164                   61           161
4 196   M     74    175                   71           175
5 197   M     83    180                   80           180
6 198   M     81    175                  NA            NA
7 199   M     90    181                   91           178
8 200   M     79    177                   81           178
9 > mean(Davis$height)
10 [1] 170.565
```

# Fichier STT1000.R

Les fichiers .R contiennent des codes R.

```
1 > url = "http://freakonometrics.free.fr/STT1000.R"  
2 > download.file(url,"STT1000.R")  
3 > source("STT1000.R")
```



# Base du langage R

```
1 > sum(Davis$sex=="M")
2 [1] 88
3 > mean(Davis$height [Davis$sex=="M"])
4 [1] 178.0114
5 > mean(Davis$height [Davis$sex=="F"])
6 [1] 164.7143
```

```
1 > aggregate(Davis$height , by=list(Davis$sex) , FUN=mean)
2   Group.1      x
3 1       F 164.7143
4 2       M 178.0114
```

```
1 > tapply(Davis$height , Davis$sex , mean)
2           F          M
3 164.7143 178.0114
```

Pour un cours de R, parcourir le livre d'Ewen Gallic.

# Base du langage R

Par exemple, on peut programmer une descente de gradient,

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} \text{ à partir d'une valeur initiale } x_0$$

```
1 > v = c(0.89367, -1.04729, 1.97133, -0.38363, 1.65414)
2 > f = function(x) sum((v-x)^2)
3 > df = function(x) -2*sum((v-x))
4 > d2f = function(x) 2*length(v)
5 > x = rep(2,100)
6 > for(i in 2:100) x[i]=x[i-1]-df(x[i-1])/d2f(x[i-1])
7 > x[100]
8 [1] 0.617644
```

ou

```
1 > df = function(x, h=1e-5) (f(x+h)-f(x))/h
2 > d2f = function(x, h=1e-5) (df(x+h)-df(x))/h
```

les fonctions, les boucles... un peu hors sujet ici, mais indispensable

## Et excel ?

A priori excel peut servir pour des statistiques descriptives, mais ce n'est pas ce qu'on va faire dans ce cours !

Avec un peu de courage, il doit être possible de refaire l'intégralité de ce que l'on verra dans un classeur excel...

Par exemple, on peut "programmer" une descente de gradient,

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} \text{ à partir d'une valeur initiale } x_0$$

	E5					
1						
2						
3	0,89367	x0		df		df2
4	-1,04729	x1	0,6176	13,82356	-8,88178E-16	10
5	1,97133	x2	0,6176		0	10
6	-0,38363					
7	1,65414					

## Et excel ?

Mais les formules sont cachées dans les cellules, ce qui rend un classeur excel difficilement auditabile

	A	B	C	D	E	F	G
1							
2					df		df2
3		0,89367	x0	2	2,21266	10	
4		-1,04729	x1	1,7787	5,652048	10	
5		1,97133	x2	1,2135	-1,5156016	10	
6		-0,38363	x3	1,3651	3,49743872	10	
7		1,65414	x4	1,0153	-1,277589024	10	
8			x5	1,1431	#VALEUR!		10

cf <https://www.wired.co.uk/article/spreadsheet-excel-errors>  
ou <https://www.bbc.com/news/magazine-22213219>

*“teaching Excel in actuarial undergraduate programs is important, but only in courses dedicated to operational risk”*