

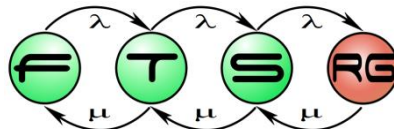
# R

„Big Data” elemzési módszerek

Kocsis Imre

[ikocsis@mit.bme.hu](mailto:ikocsis@mit.bme.hu)

2015.09.30.



# Adatelemzés (a számítógépig)

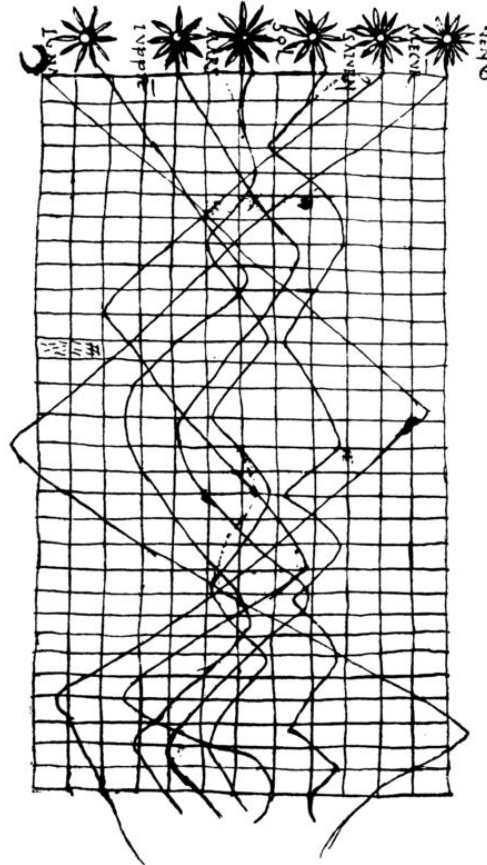
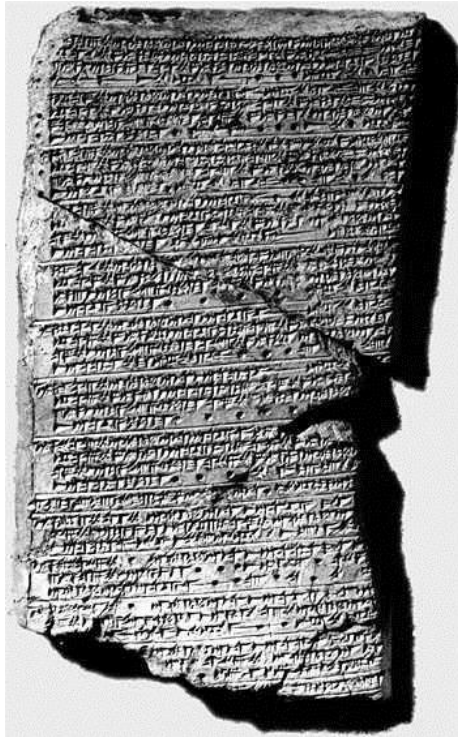
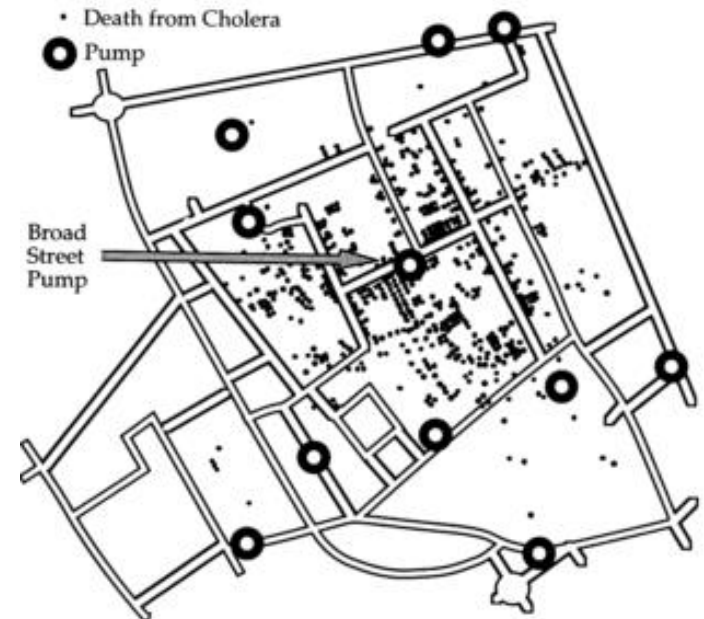


Table 1  
Cholera deaths in Golden Square, Soho from August 31 to September 11, 1854 (adapted from Snow, 1855a, p. 49)

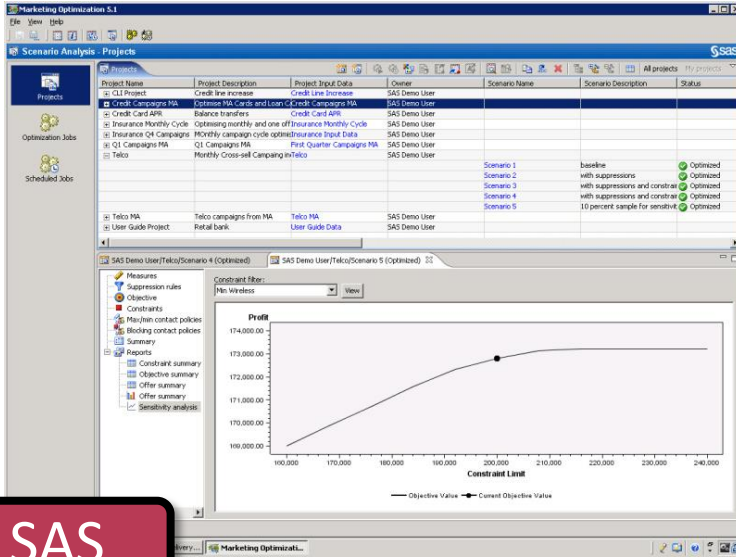
Date	No. of fatal attacks	Deaths
August 31	56	3
September 1	143	70
September 2	116	127
September 3	54	76
September 4	46	71
September 5	36	45
September 6	20	37
September 7	28	32
September 8	12	30
September 9	11	24
September 10	5	18
September 11	5	15



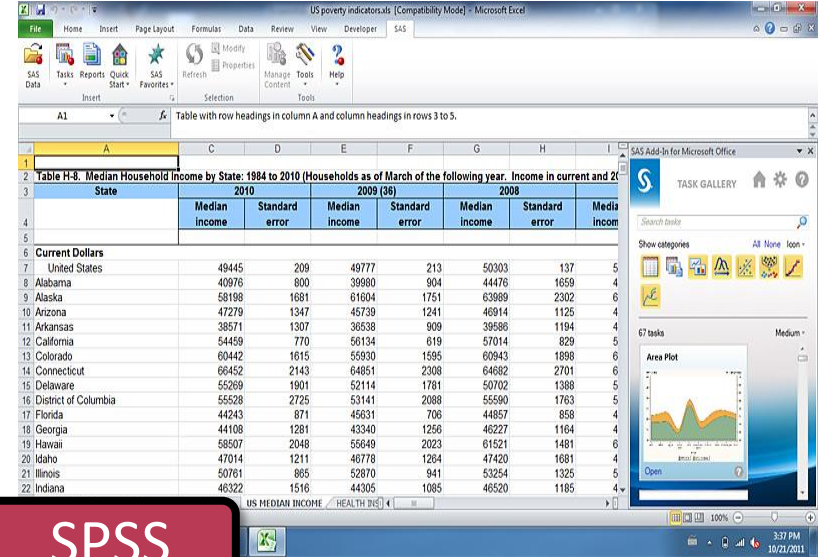
# Mi *nem* statisztikai eszköz/csomag?

- Táblázatkezelő
  - Lásd pl. [4]
- Adatbáziskezelő
  - SQL
- Saját C/FORTRAN/Perl/Java...
  - EDA...?
  - Stat. függvények?
- Úgy értve, hogy klasszikusan
  - Mindhárom területen változik
- + adatelemzés != statisztika

# Mi az, ami igen



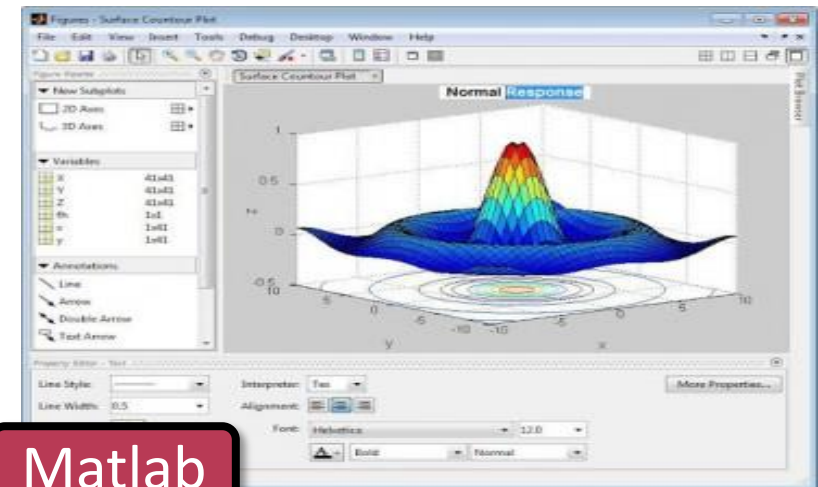
SAS



SPSS

```
1) meetup.R
2) Source on Save
3) Run
4) Source
5) Import Dataset
6) Workspace
7) Files
8) Plots
9) Packages
10) Help
11) Data
12) 44593 obs. of 9 variables
13) faulty
14) 153 obs. of 9 variables
15) ok
16) 44440 obs. of 9 variables
17) Values
18) fts
19) iset[9]
20) myfilepath
21) "C:/drive/ducroot/FTSRG/SessionStore/2013_06_16_thiny_multicheck/"
22) oks
23) iset[9]
24) Console
25) C:/drive/ducroot/FTSRG/SessionStore/2013_06_16_thiny_multicheck/
26) > ls()
27) [1] "dat" "faulty" "fts"
28) [4] "myfilepath" "ok" "oks"
29) >
```

R



Matlab

+ wikipedia [5]

# Néhány általános jellemző

- Saját szkriptnyelv
  - Interaktív futtatással is
- *Validált* stat. eljárások széles köre
  - As in: „clinical trial data for FDA submissions”
- „Workspace” modell
  - Jellemzően in-memory (vs. „out-of-memory” elemzés)
- Erős vizualizációs képességek
- Kapcsolódó funkciók
  - jelentések, adatbázis-kapcsolat, GUI-szkriptelés, webalkalmazások, munkafolyamatok, etc.
- Gyökerek: 70-es évektől →...
  - SAS Institute: 1976, az egyetemmel szemben
  - Szoftvertechnológiailag erősen látszik; az új generáció már más
  - <http://julialang.org/blog/2012/02/why-we-created-julia/>

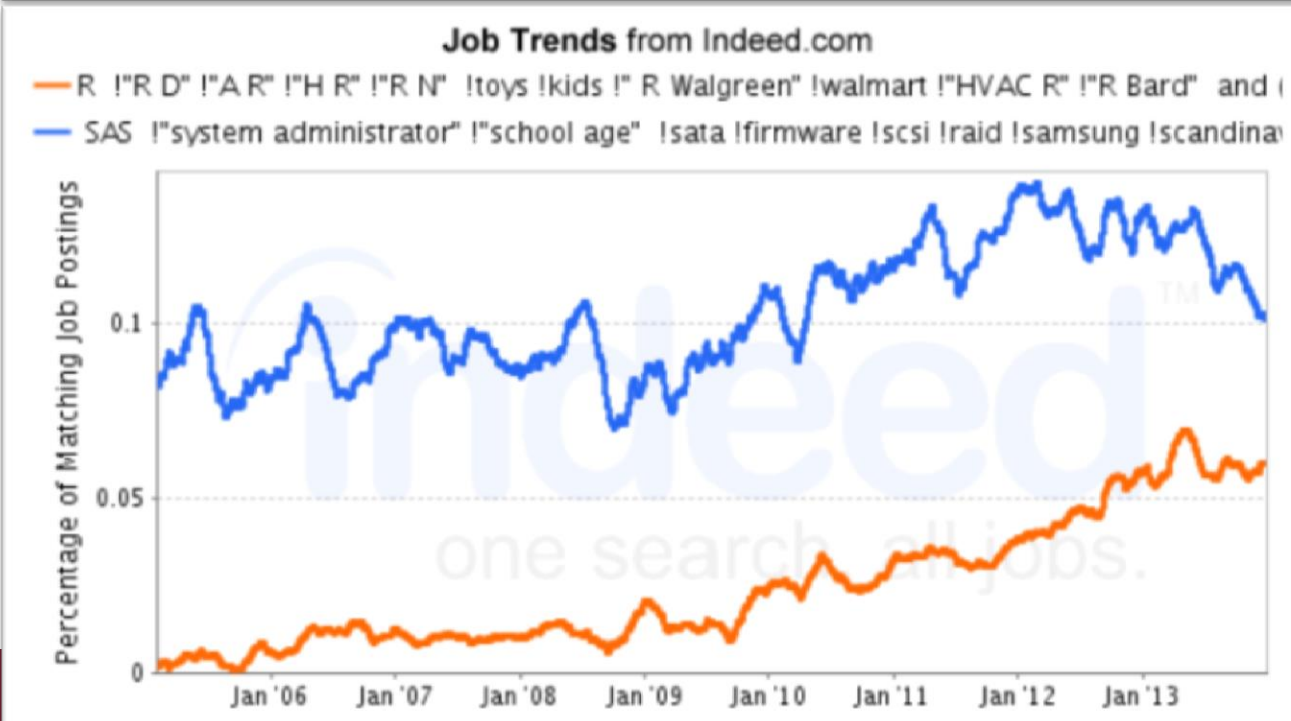
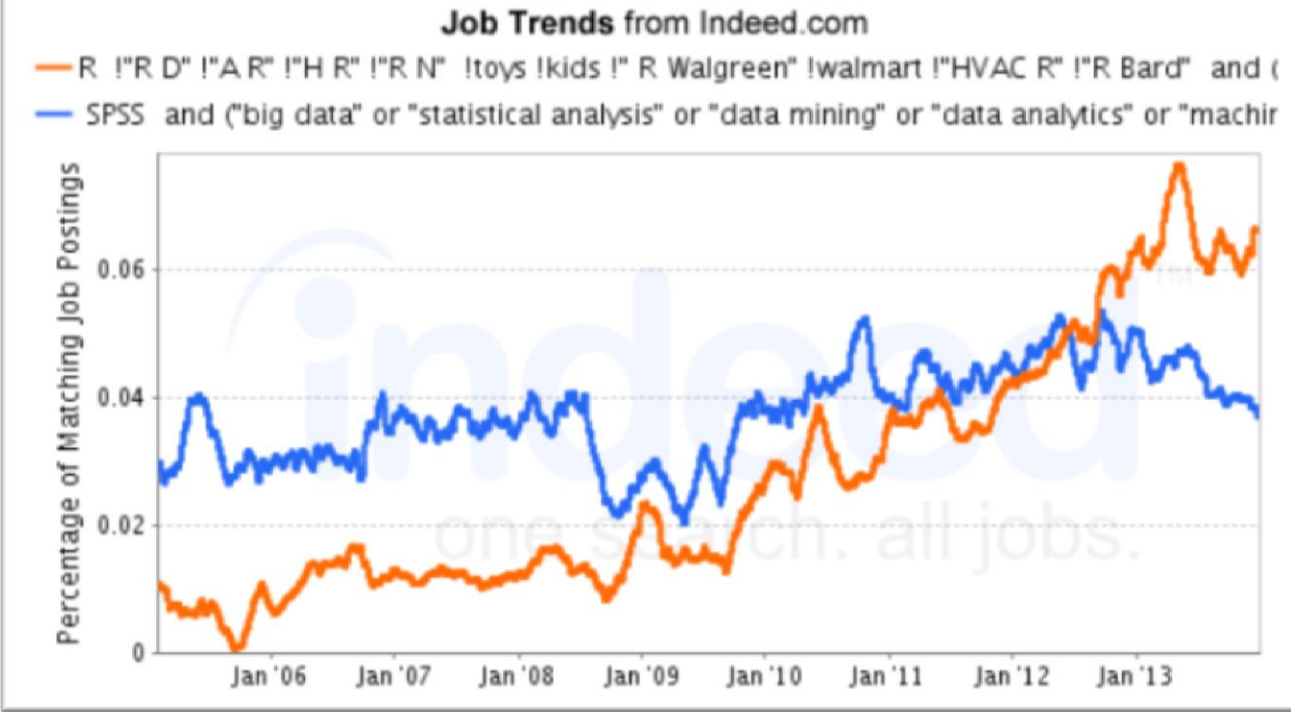
# R

- Az S nyelv „GNU verziója”
  - Környezet és nyelv egyben
- Statisztikai számítások és grafika
- Nem csak ingyenes; *nyílt* is
- Hatékonyság: „kihívás” C/C++/FORTRAN-ba
- Egyre inkább „lingua franca”, ha adatot kell elemezni
  - + Python





# Miért R? (r4stats.com)



# Miért R? (r4stats.com)

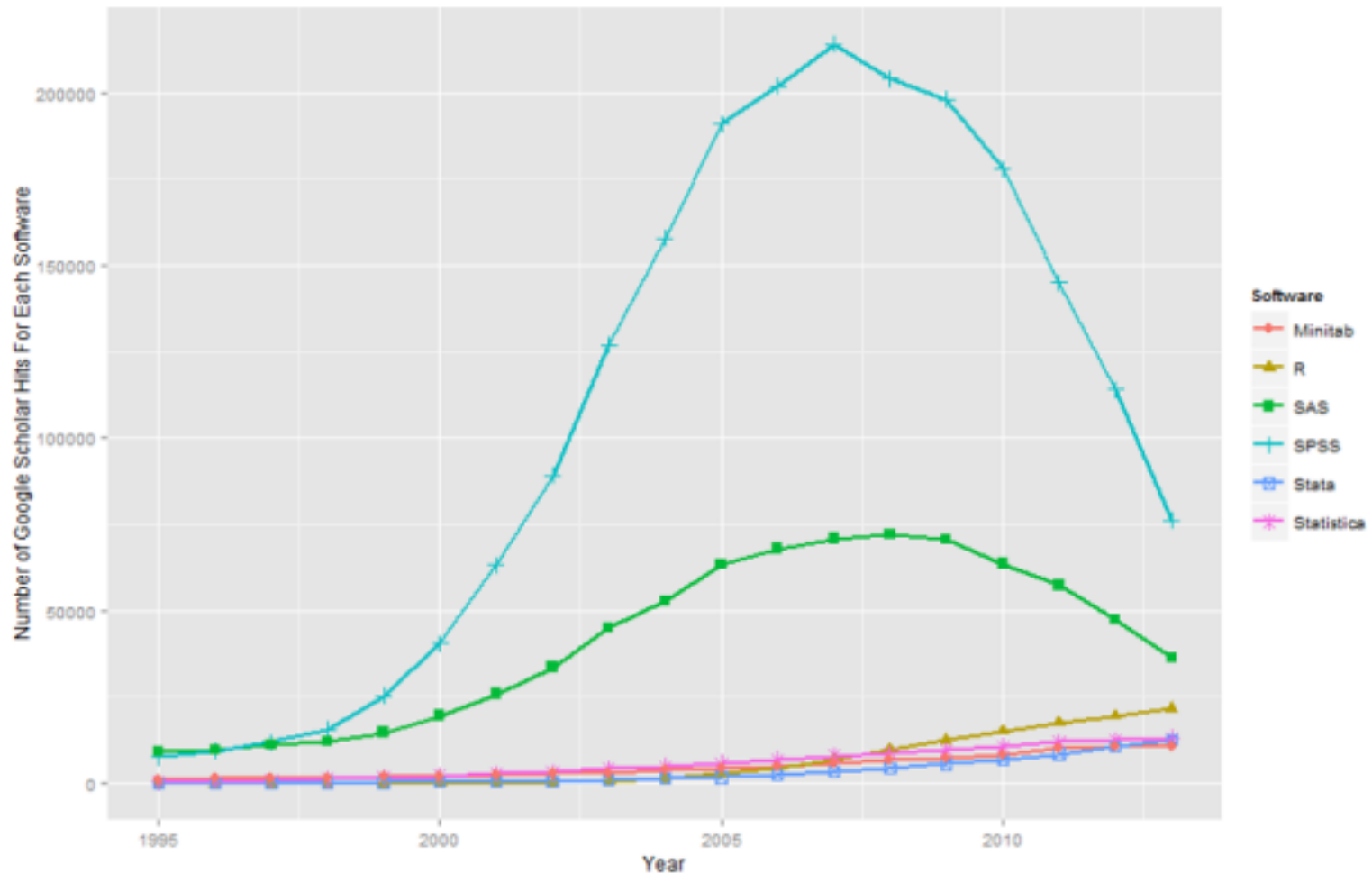


Figure 2b. Number of scholarly articles found for the top five classic statistics packages.



# Miért R? (r4stats.com)

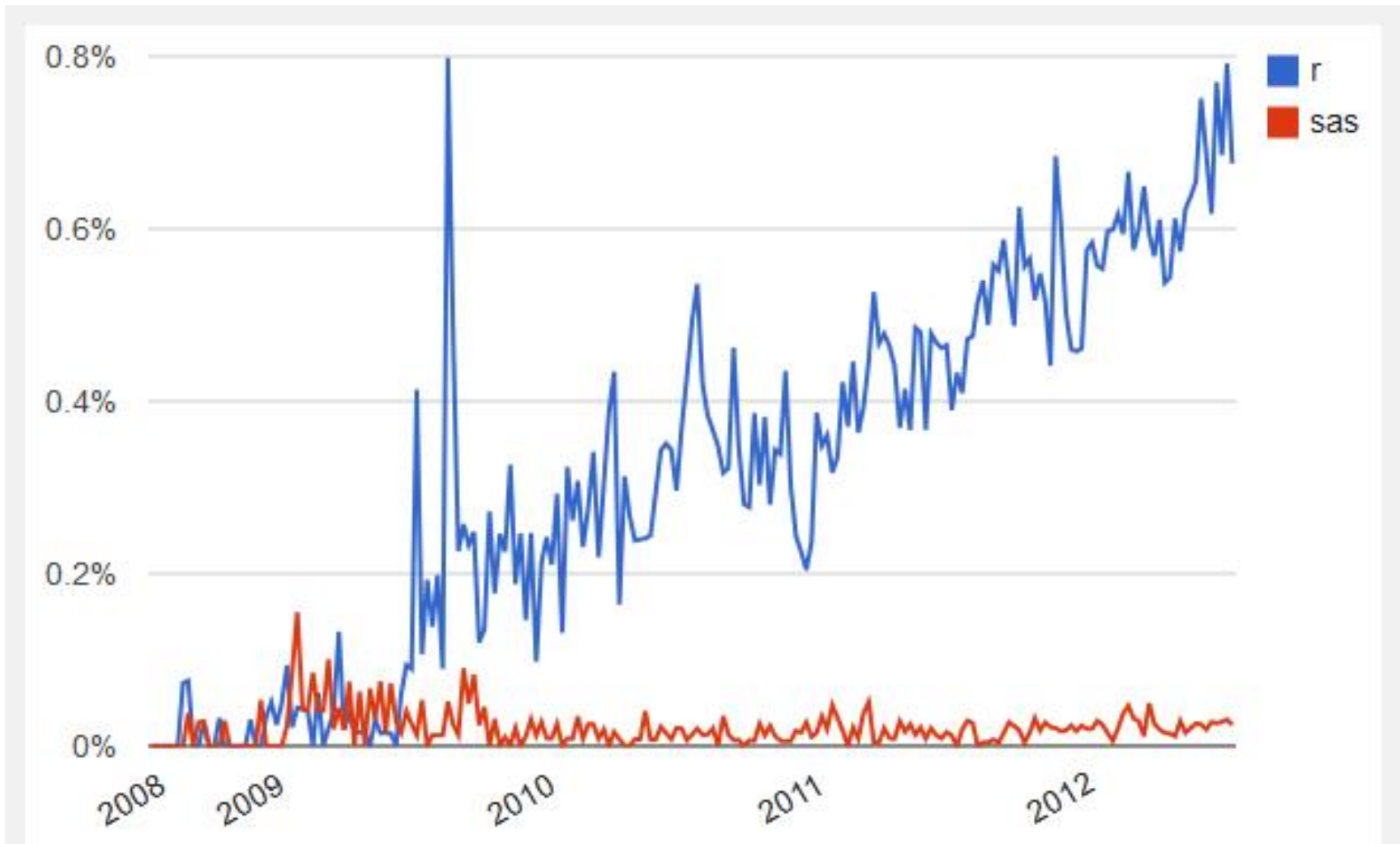


Figure 1c. Number of R- or SAS-related posts to Stack Overflow by week.

Forrás: [1]

# Miért R? (r4stats.com)

- Mert HF

Forrás: [1]



## The Comprehensive R Archive Network

### Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

### Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2013-09-25, Frisbee Sailing) [R-3.0.2.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features and bug fixes](#) before filing corresponding feature requests or bug reports.

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](#)

<a href="#">Bayesian</a>	Bayesian Inference
<a href="#">ChemPhys</a>	Chemometrics and Computational Physics
<a href="#">ClinicalTrials</a>	Clinical Trial Design, Monitoring, and Analysis
<a href="#">Cluster</a>	Cluster Analysis & Finite Mixture Models
<a href="#">DifferentialEquations</a>	Differential Equations
<a href="#">Distributions</a>	Probability Distributions
<a href="#">Econometrics</a>	Computational Econometrics
<a href="#">Environmetrics</a>	Analysis of Ecological and Environmental Data
<a href="#">ExperimentalDesign</a>	Design of Experiments (DoE) & Analysis of Experimental Data
<a href="#">Finance</a>	Empirical Finance
<a href="#">Genetics</a>	Statistical Genetics
<a href="#">Graphics</a>	Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization
<a href="#">HighPerformanceComputing</a>	High-Performance and Parallel Computing with R
<a href="#">MachineLearning</a>	Machine Learning & Statistical Learning
<a href="#">MedicalImaging</a>	Medical Image Analysis
<a href="#">MetaAnalysis</a>	Meta-Analysis
<a href="#">Multivariate</a>	Multivariate Statistics
<a href="#">NaturalLanguageProcessing</a>	Natural Language Processing
<a href="#">NumericalMathematics</a>	Numerical Mathematics
<a href="#">OfficialStatistics</a>	Official Statistics & Survey Methodology
<a href="#">Optimization</a>	Optimization and Mathematical Programming
<a href="#">Pharmacokinetics</a>	Analysis of Pharmacokinetic Data
<a href="#">Phylogenetics</a>	Phylogenetics, Especially Comparative Methods
<a href="#">Psychometrics</a>	Psychometric Models and Methods
<a href="#">ReproducibleResearch</a>	Reproducible Research
<a href="#">Robust</a>	Robust Statistical Methods
<a href="#">SocialSciences</a>	Statistics for the Social Sciences
<a href="#">Spatial</a>	Analysis of Spatial Data
<a href="#">SpatioTemporal</a>	Handling and Analyzing Spatio-Temporal Data
<a href="#">Survival</a>	Survival Analysis
<a href="#">TimeSeries</a>	Time Series Analysis
<a href="#">WebTechnologies</a>	Web Technologies and Services
<a href="#">gR</a>	gRaphical Models in R

### Available Packages

Currently, the CRAN package repository features 5898 available packages.

+ GitHub, BioC, R-Forge, saját, ...

# R konzol

R i386 2.15.3

R x64 2.15.3

```
R version 2.15.3 (2013-03-01) -- "Security Blanket"
Copyright (C) 2013 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: x86_64-w64-mingw32/x64 (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.

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.

> |
```

...

# RStudio

The image shows the RStudio interface with three callouts pointing to different parts of the software:

- Parancsállományok** (Script files): Points to the source editor window showing R code.
- Interaktív konzol** (Interactive console): Points to the console window showing the output of the `ls()` command.
- „workspace”** (Workspace): Points to the Environment pane showing the current workspace objects.

```
73 ihist(dat$RPT)
74
75 datiset <- iset
76 factors <- sapp
77 numerics <- sapp(numerics, func=fun(x){datiset[[x]]})
78
79 sapply(factors, ibar)
80 sapply(numerics, ihist)
81
82 fts <- iset.new()
83 ihist(fts$RT, ti
84
85
```

```
> ls()
[1] "dat"      "faulty"   "fts"
[4] "myfilepath" "ok"      "oks"
>
```

Data	
dat	44593 obs. of 9 variables
faulty	153 obs. of 9 variables
ok	44440 obs. of 9 variables
Values	
fts	iset[9]
myfilepath	"C:/gdrive/docroot/FTSRG/s
oks	iset[9]

```
?smooth.spline
smooth.spline(iris$Sepal.Length)
?loess
library("bigvis", lib.loc="C:/User.
ls()
```



# Ismerkedés az R-rel

- Interaktív bevezetés az R nyelvbe és környezetbe példákon keresztül
- Rintro.R
- Induláshoz javasolt:
  - FTSRG tech cheat sheet [6]
  - Magyarul: [2] és [3]
- N.B.: nem kell hozzá informatikusnak lenni
  - Előny és hátrány is

# typeof

## 1. táblázat. Fontosabb typeof visszatérési értékek

érték	jelentése
NULL	Null
symbol	változó neve
closure	függvény
logical	logikai értékekből álló vektor
integer	egész számokból álló vektor
double	lebegőpontos számokból álló vektor
complex	komplex adatokból álló vektor
character	karaktervektor
list	lista
raw	bináris vektor

Forrás: [2]



**Find**  
a Meetup Group

**Start**  
a Meetup Group



# Budapest Users of R Network



Home

Members

Sponsors

Photos

Discussions

More

 My profile



## Budapest, Hungary

Founded Aug 11, 2013

About us...

useRs 216

Group reviews 2

### > print('Hello, Hungary!')

+ SUGGEST A NEW MEETUP

Upcoming 1

Suggested 0

Past

Calendar

### R integráció kereskedelmi és egyéb termékekben

BME I. épület, I.B 019. terem (földszinten a porta mögött)

BME I épület, Budapest, Magyar tudósok körútja, Hungary, Budapest (map)



Mon Oct 13  
6:00 PM

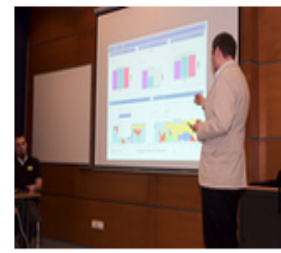
I'M GOING

39 going

1 comment

This will be a Hungarian speaking event,

### What's new



# R for big data (abridged)

Visualization ⊕

C API  
Rcpp ⊕ Efficiency

parallel  
Rmpi ⊕  
foreach ⊕ Parallel

⊕ Teradata  
⊕ Oracle  
Netezza  
⊕ SAP HANA  
in-database analytics

⊕ GPU

⊕ Rserve  
⊕ Python  
⊕ Perl  
⊕ Web  
⊕ SAS  
As backend

⊕ Rgraphviz  
Rcpp  
rJava  
rPython  
Glue

Data formats

Flat text ⊕  
HDFS ⊕  
SQL ⊖  
  sqldf  
  RODBC  
  Rmysql  
  RJDBC  
  ROracle  
NoSQL ⊖  
  MongoDB ⊕  
  CouchDB ⊕  
JSON ⊕  
XML ⊕  
HBase ⊕

Large & out-of-memory data

ff  
bigmemory  
biglm  
biglars  
bigrf

Hadoop

RHIPE  
rmr  
HadoopStreaming 📄  
Revolution Analytics ⊕

# Hivatkozások

- [1] <http://r4stats.com/articles/popularity/>
- [2] <http://cran.r-project.org/doc/contrib/Solymosi-Rjegyzet.pdf>
- [3] <http://www.inf.unideb.hu/~jeszy/R/>
- [4] <http://people.umass.edu/evagold/excel.html>
- [5] [http://en.wikipedia.org/wiki/List\\_of\\_statistical\\_packages](http://en.wikipedia.org/wiki/List_of_statistical_packages)
- [6] <https://github.com/FTSRG/technology-cheat-sheets/wiki/R-programming-language>