

Utilizarea programului R pentru analiză statistică



Daniel-Corneliu Leucuța, MD, PhD, lector
Departamentul de Informatică Medicală și Biostatistică
Universitatea de Medicină și Farmacie "Iuliu Hațieganu", Cluj-Napoca
dleucuta@umfcluj.ro

Conținut

R utilizare, popularitate

Disponibilitate, instalare

Interfață grafică cu exemple în R Commander

- Lucrul cu datele (importare, gestiunea setului de date și a variabilelor)
- Statistică de bază (stastică descriptivă, Teste de normalitate, proporții, medii, varianțe, neparametrice, corelații)
- Statistică avansată (modele, diagnostica prezumpțiilor),
- Grafică,
- Afișarea rezultatelor

R/ R Commander pentru a preda statistica

Utilizarea R în linie de comandă

Pachete utile

Găsirea documentației, găsirea de ajutor cu funcții și pachete

Ce este R

R este o implementare a limbajului S

Oferă un număr imens de tehnici de analiză statistică și reprezentare grafică

(11,000 packages (July 2017) disponibile pe Comprehensive R Archive Network (CRAN), Bioconductor, Omegahat, GitHub, ...)

Organizat, flexibil, extensibil, coerent

Programare orientată obiect

Limbaj interpretat (o parte din pachete scrise în C/C++ pentru creșterea vitezei)

Gratuit

Avantaje și dezavantaje

Avantaje

- Gratuit
- Zeci de mii de pachete cu funcții statistice
- Ușor de extins
- Capabilități grafice
- În trend cu machine learning

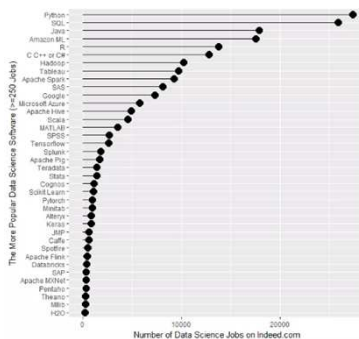
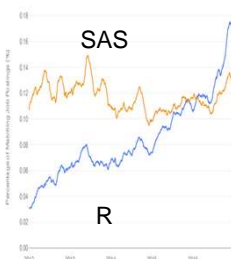
Dezavantaje

- Securitate deficitară
- Gestiunea memoriei deficitară

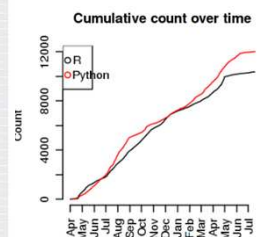
Deși în linie de comandă (majoritar), nu e dificil de folosit de persoane fără experiență în programare

Popularitate

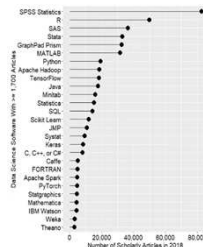
Joburi în data science



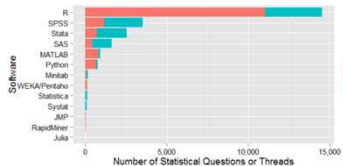
Kaggle data science contests



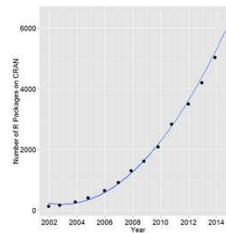
Articole științifice



Întrebări pe site-uri



Pachete în R



<http://r4stats.com/articles/popularity/> The Popularity of Data Science Software. by Robert A. Muenchen

Disponibilitate

Disponibil pe numeroase platforme:

- Windows
- Linux
- OS X
- Android

Ce firme cunoscute folosesc R

Microsoft Google facebook Tweeter

IBM

Companies that use R for Analytics

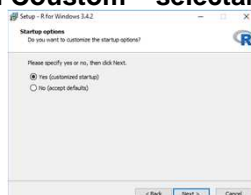


Instalare

Pe Windows – instalare ușoară

<http://www.r-project.org/>

- pe prima pagină - link download R
- În pagina deschisă (CRAN Mirrors) selectat un server
- În noua pagină se selectează link-ul pentru Windows.
- În nouă pagină, R for Windows, unde se urmează legătura numită base (binaries for base distribution).
- În nouă pagină, intitulată *R-numărul versiunii for Windows (32/64 bit build)* (ex. *R-3.6.1 for Windows*). Se urmează link-ul din susul paginii *Download R numărul versiunii for Windows*.
- Instalare, cu Coustom – selectare SDI, în rest next, next, finish.



Instalare

Pe Windows

- Există și versiunea de la Microsoft – Microsoft R Open

Pe Linux

<http://www.jason-french.com/blog/2013/03/11/installing-r-in-linux/>

Pe Mac - OSx

<https://cran.r-project.org/bin/macosx/RMacOSX-FAQ.html>

Pe Android

<https://www.r-bloggers.com/install-r-in-android-via-gnuroot-no-root-required/>

Interfețe grafice

Covârșitor funcțiile din R sunt accesibile prin linie de comandă

Există câteva interfețe grafice:

- Mediu de dezvoltare grafic integrat: Rstudio Desktop (gratuit) (desktop/server – cu plată – suport/cloud)
- Editoare R
 - Notepad ++ (Windows)
 - Tinn-R (Windows)
 - Eclipse (numeroase platforme)
 - Geany (Linux)
- Interfețe Point and click (Windows, Linux, Mac OS):
 - [R Commander](#) – cea mai cuprinzătoare și logic structurată
 - [RKWard](#) – ca un SPSS, dar relativ puține funcții
 - [Rattle GUI](#) - R Analytical Tool To Learn Easily – pentru data mining (cluster analysis, association analysis, decision trees, random forests, boosting, support vector machines, model performance evaluation, exploring data)

R Commander

Interfață point and click

Frumos structurată

Probabil cea mai cuprinzătoare (extinsă prin plug-in-uri >40, relativ ușor de construit)

```
[1] "Rcmdr"
[5] "RcmdrPlugin.BiclustGUI"
[9] "RcmdrPlugin.doex"
[13] "RcmdrPlugin.epack"
[17] "RcmdrPlugin.FuzzyClust"
[21] "RcmdrPlugin.KMggplot2"
[25] "RcmdrPlugin.MPASTate"
[29] "RcmdrPlugin.plotByGroup"
[33] "RcmdrPlugin.RMTCJags"
[37] "RcmdrPlugin.SLC"
[41] "RcmdrPlugin.survival"
[45] "RcmdrPlugin.UCA"

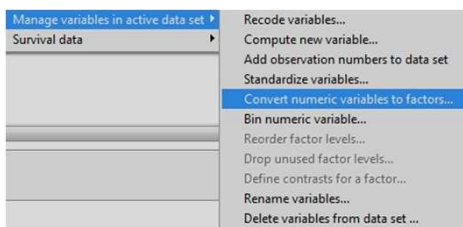
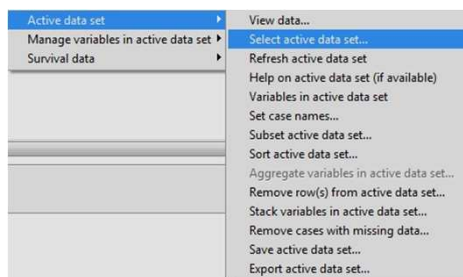
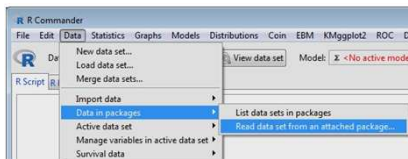
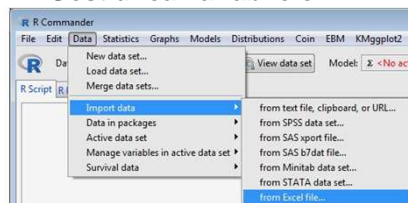
"RcmdrMisc"
"RcmdrPlugin.coin"
"RcmdrPlugin.EACSPIR"
"RcmdrPlugin.Export"
"RcmdrPlugin.GWRM"
"RcmdrPlugin.lfstst"
"RcmdrPlugin.NMBU"
"RcmdrPlugin.pointG"
"RcmdrPlugin.BOC"
"RcmdrPlugin.SM"
"RcmdrPlugin.sutteForecastR"

"RcmdrPlugin.aRnova"
"RcmdrPlugin.depthTools"
"RcmdrPlugin.EBM"
"RcmdrPlugin.EZR"
"RcmdrPlugin.HH"
"RcmdrPlugin.MA"
"RcmdrPlugin.orloca"
"RcmdrPlugin.qual"
"RcmdrPlugin.sampling"
"RcmdrPlugin.sos"
"RcmdrPlugin.TeachingDemos"

"RcmdrPlugin.BCA"
"RcmdrPlugin.DoE"
"RcmdrPlugin.EcoVirtual"
"RcmdrPlugin.FactoMineR"
"RcmdrPlugin.IPSUR"
"RcmdrPlugin.mosaic"
"RcmdrPlugin.PcaRobust"
"RcmdrPlugin.RiskDemo"
"RcmdrPlugin.SCDA"
"RcmdrPlugin.steepness"
"RcmdrPlugin.temis"
```

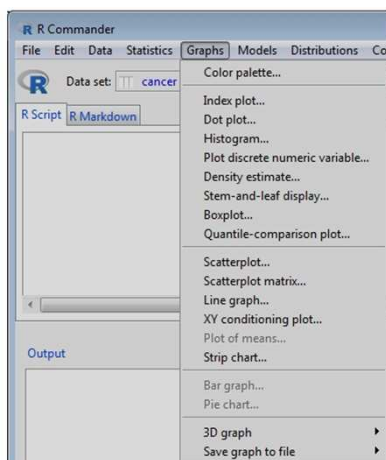
R Commander

- Lucrul cu datele
- Importare
- Gestiunea setului de date
- Gestiunea variabilelor



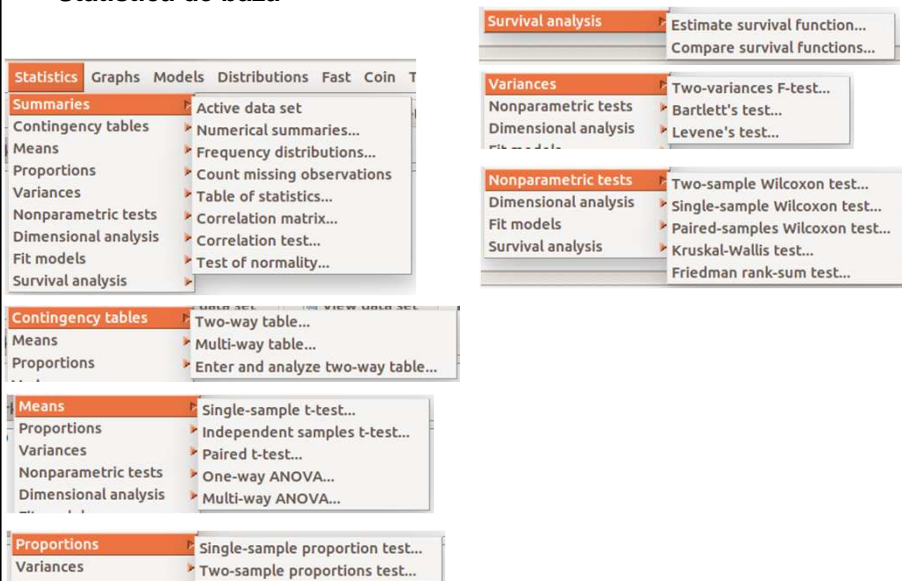
R Commander

- Grafică



R Commander

- Statistică de bază



R Commander

- Statistică avansată
 - (modele, diagnostica prezumpțiilor)

The screenshot displays the R Commander interface with several menus open:

- Dimensional analysis:** Scale reliability..., Principal-components analysis..., Factor analysis..., Confirmatory factor analysis..., Cluster analysis
- Fit models:** Linear regression..., Linear model..., Generalized linear model..., Multinomial logit model..., Ordinal regression model..., Cox regression model..., Parametric survival model...
- Hypothesis tests:** ANOVA table..., Compare two models..., Linear hypothesis...
- Numerical diagnostics:** Variance-inflation factors, Breusch-Pagan test for heteroscedasticity..., Durbin-Watson test for autocorrelation..., RESET test for nonlinearity..., Bonferroni outlier test, Test proportional hazards
- Models:** Select active model..., Summarize model, Compare model coefficients..., Add observation statistics to data..., Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), Stepwise model selection..., Subset model selection..., Confidence intervals..., Bootstrap confidence intervals..., Delta method confidence interval..., Hypothesis tests, Numerical diagnostics, Graphs
- Graphs:** Basic diagnostic plots, Residual quantile-comparison plot..., Component-residual plots, Added-variable plots, Influence plot, Effect plots, Cox-model survival function..., Plot terms in Cox model, Plot survival-regression dfbetas, Plot survival-regression dfbeta, Plot null Martingale residuals, Cox-model partial-residual plots

R Commander

Afișarea rezultatelor

- Output (text)
- Fereastra de grafice
- R Markdown
 - Html
 - Doc (dacă e instalat pandoc)
 - PDF (dacă e instalat LaTeX)

The screenshot shows the R Commander interface with the R Markdown editor and the Output window. The R Markdown code is as follows:

```

<!-- R Commander Markdown Template -->
Replace with Main Title
-----
### Your Name
### "r" as.character(Sys.Date())
### (r echo=FALSE)
  
```

The Output window displays the following text:

```

quantile.Survfit.quantiles<(25, 50, 75)>
  Squantile      25  50  75
  Tip:celule mici 29  51  99
  Tip:scuamos    33 118 357
  slower         25  50  75
  Tip:celule mici 13 25  61
  Tip:scuamos    11 82 228
  Supper        25  50  75
  Tip:celule mici 27  61 153
  Tip:scuamos    112 314 991
  
```

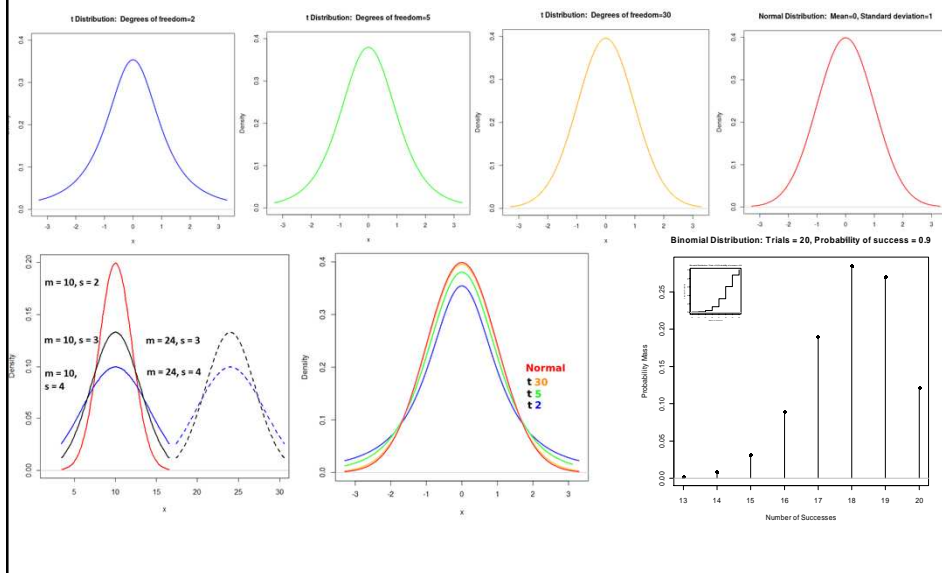
Below the text, a survival plot is shown, comparing the survival curves for 'celule mici' (solid line) and 'scuamos' (dashed line). The x-axis represents time (0 to 1000) and the y-axis represents survival probability (0.0 to 1.0).

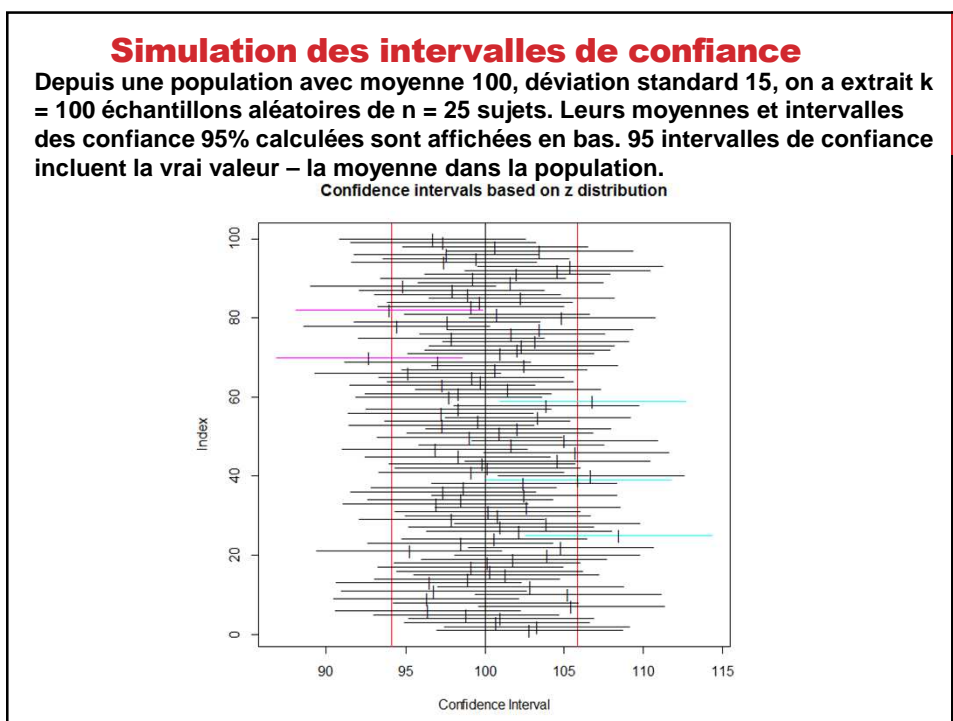
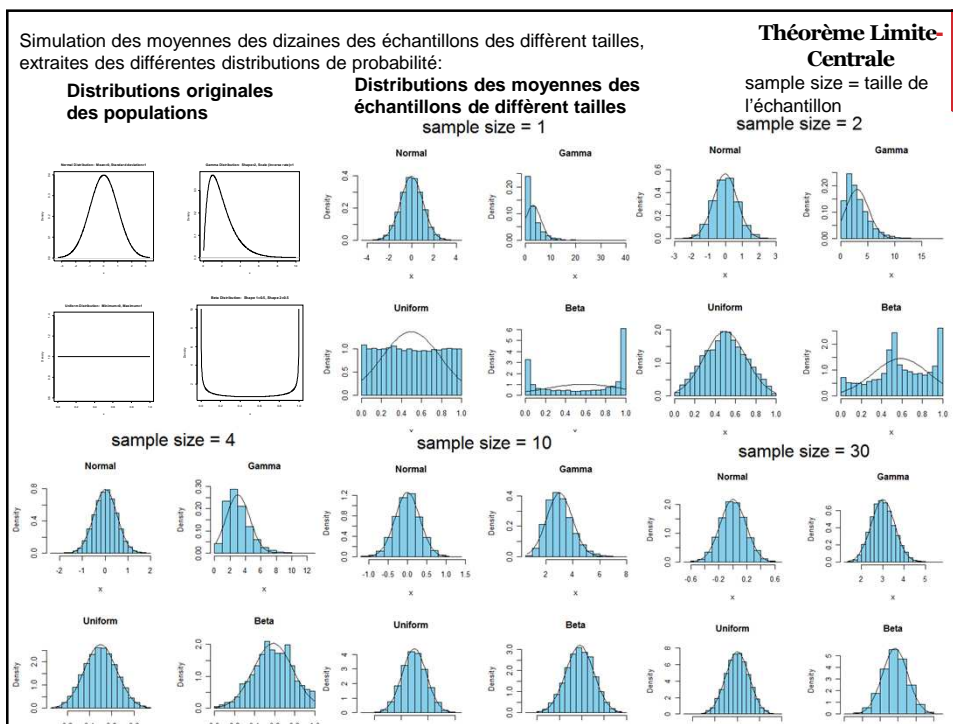
R/R Commander pentru a preda statistica

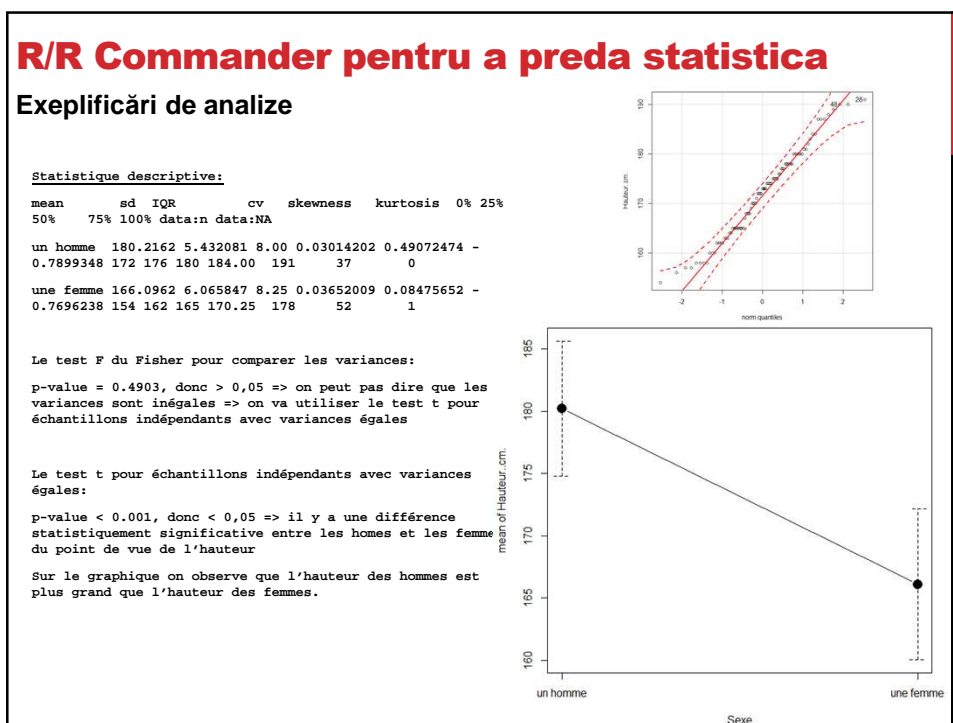
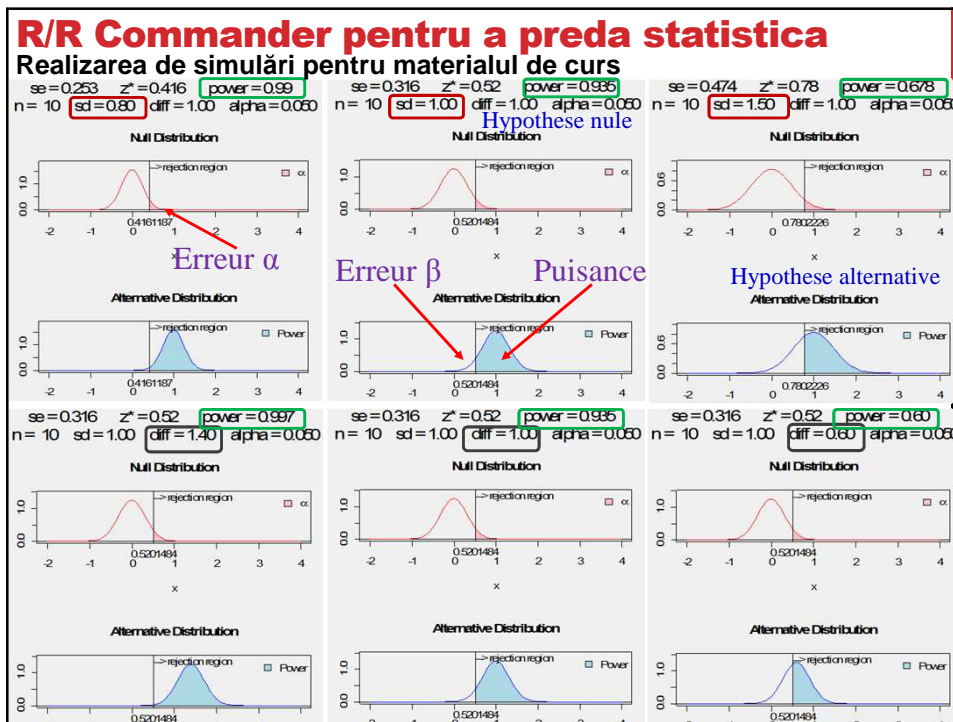
- Realizarea de simulări pentru materialul de curs
 - Distribuții de probabilitate – meniul Distributions ...
 - Plug-in RcmdrPlugin.TeachingDemos
- Exeplificări de analize
 - R Commander ...

R/R Commander pentru a preda statistica

Realizarea de simulări pentru materialul de curs







Utilizarea R în linie de comandă

Nu e așa de dificil cum am putea crede, nu trebuie să fim programatori

Esența pentru începători

- Se importă baza de date cu Rcommander
- Variabilele se identifică prin: nume.dataset\$nume.variabilă
- Se caută documentația pentru funcția de interes
- Se înlocuiesc/introduc numele variabilelor
- Se înlocuiesc/introduc numele setului de date
- Se precizează opțiunile

Ex: testul student pentru eșantioane independente:

- `t.test(dataset$age~dataset$gender, alternative='two.sided', conf.level=.95, var.equal=FALSE)`

sau

- `t.test(age~gender, alternative='two.sided', conf.level=.95, var.equal=FALSE, data=dataset)`

Utilizarea R în linie de comandă

Variabile

```
> x <- 1
> x
[1] 1
> x <- "abc"
> x
[1] "abc"
> x <- seq(from=1,to=5,by=.5)
> x
[1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0
> x <- rep("A",5)
> x
[1] "A" "A" "A" "A" "A"
> x <- TRUE
> x <- c("A", "A", "b")
```

Utilizarea R în linie de comandă

Datele lipsă

NA - ideal de înlocuit manual spațiile goale din fișierul Excel cu NA

Comentarii – semnul

```
x <- 2x
# x <- 1
x
```

Utilizarea R în linie de comandă

Matrice

```
x = matrix(
  c(1, 2, 3, 4, 5, 6),
  nrow=3,
  ncol=2)
x
```

```
y = matrix(
  c(3, 2, 1),
  nrow=3,
  ncol=1)
```

```
x <- cbind(x,y) - adauga o coloana
```

Dataframe

Utilizarea R în linie de comandă

Funcții definite de utilizator

```
fncDivideBy2 <- function(x) {  
    return(x/2)  
}  
x <- 2  
print(fncDivideBy2(x))  
[1] 1  
  
fncAxB <- function(a,b) {  
    return(a*b)  
}  
print(fncAxB(2,3))  
print(fncAxB(a=2,b=3))  
  
[1] 6  
[1] 6
```

Utilizarea R în linie de comandă

Programare în R – bucla for

```
x <- c(5, 4, 3, 2, 1)  
for(i in 1:length(x)) {  
    print(x[i]^2)  
}  
print(i)  
  
[1] 25  
[1] 16  
[1] 9  
[1] 4  
[1] 1  
  
> print(i)  
[1] 5
```

Utilizarea R în linie de comandă

Programare în R – bucla while

```
i <- 1
while(i<5) {
  print(i^2)
  i <- i+1
}
print(i)
```

Utilizarea R în linie de comandă

Programare în R – execuție condiționată

```
x <- 1
if (x==1) {
  print("x=1")
} else {
  print("x!=1")
}

x <- 1
y <- ifelse(x==1,y=1,y=0)
```

Pachete utile

Sunt nenumărate pachete utile a fi utilizate în R, depinde mult de domeniul în care lucrăm

Există liste de pachete pe site-ul programului:

<https://cran.r-project.org/web/views/>

Unele pachete utile în statistica medicală ar putea fi:

Grapher - pentru realizare de grafice cu calitate înaltă pentru a publica

LogisticDx - pentru analiza diagnostică a regresiiilor logistice

Google ;)

Găsirea documentației, găsirea de ajutor cu funcții și pachete

Google R + întrebare, sau nume pachet R + întrebare, sau nume funcție R + întrebare

[Stackoverflow.com](https://stackoverflow.com) - site cu întrebări și răspunsuri (include și R)

în R

```
help(t.test)
```

```
help("t.test")
```

```
?t.test
```

```
? "t.test,"
```

Utilizarea de "vignete" - PDF-uri cu tutoriale, materiale didactice, de clarificare

browseVignettes() - pentru a afla toate vignete-le disponibile în pachetele instalate

Utilizarea demonstrațiilor

```
demo() pentru a afla toate demonstrațiile
```