

INIȚIERE în R

Pentru persoane cu pregătire matematică

Maria Miroiu

Viorel Petrehuș

Gheorghiță Zbăganu

Introducere

Mediul de programare “R” face parte dintr-un “proiect” mai mare de a produce software gratuit. Se poate spune că el este un dialect al altui limbaj, numit „S”, creat în 1976 la laboratoarele Bell pentru a se ușura calculele statistice. La început s-a folosit o sintaxă similară Fortran-ului, dar în 1988 a fost rescris în C. Spre deosebire de „R”, „S” nu este un software gratuit. Prima variantă de „R” a fost creată în Noua Zeelandă, 1991, de Ross Ihaka și Robert Gentleman ¹. În 1993 s-a făcut lansarea oficială a proiectului, conceput la început ca un software statistic non – gratuit, așa cum mai erau și altele - SAS, SPSS, STATISTICA. În 1995 autorii s-au hotărât să îl facă un software gratuit. Pentru a controla nucleul de bază al limbajului, s-a creat în 1997 așa numitul „R Core Group”, care controlează codul sursă. În 2000 s-a lansat versiunea R.1.0.0 iar în 2008 s-a ajuns la versiunea R.2.7.2. Noi am folosit versiunea R.2.12.2, apărută pe 25 februarie 2011. Există zeci de site-uri de unde se poate descărca sistemul R și pachete adiționale..

Părintele limbajului , Ross Ihaka scria în 1998 ²:

„În august 1993 Robert Gentleman și cu mine am plasat copii ale R-ului la Statlib și am făcut un mic grup pe s-news. Unii ne-au descărcat copiile și s-au oferit să colaboreze. Cel mai de nădejde a fost Martin Maehler de la ETH Zurich care ne-a convins să lansăm codul sursă ca un software gratuit. La început am avut ezitări, dar în iunie 1995 am convenit să punem la dispoziția tuturor codul sursă în condițiile fundației FRF ³ . ”

S-a creat astfel o impresionantă comunitate a utilizatorilor și dezvoltatorilor limbajului. Toți voluntari.

Decizia de a face R un software gratuit, continuă Ihaka *„ne-a făcut să ne propunem țeluri mai mari decât un simplu soft statistic cu care să ne putem preda cursurile mai bine. Ne-a dat acces la o masă de indivizi foarte talentați doritori să se implice în proiect. Poate unul din cele mai bune lucruri care ni s-a întâmplat a fost șansa de lucra cu asemenea oameni strălucitori ”*. Și își exprima optimismul : *„ Se pare că R a ajuns la punctul în care este destul de stabil pentru a fi folosit de utilizatori – cel puțin sub Unix. Sperăm ca în 1999 să putem lansa varianta completă R.1.0 în cadrul fundației Free Software ”*

De ce software gratuit? Ce avantaje avem?

În primul rând, firește, nu dăm bani pe el. Dar mai sunt și alte avantaje, cum arată un entuziast al R-ului, Roger. D. Peng ⁴ *“Un software gratuit vă garantează libertatea de a folosi programele lui în orice scop (libertatea 0); libertatea de a studia cum funcționează programele și de a le adapta la nevoile voastre. Accesul la codul*

¹ <http://www.biostat.jhsph.edu/~rpeng/biostat776/lecture1.pdf>

² <http://cran.r-project.org/doc/html/interface98-paper/paper.html>

³ <http://www.fsf.org/>

⁴ <http://www.biostat.jhsph.edu/~rpeng/>

sursă (libertatea 1). Libertatea de a redistribui copiile pe care le aveți oricui, așa încât să îi poți ajuta vecinul (libertatea 2); libertatea de a distribui oricui îmbunătățirile făcute de tine”.

Sigur că există și dezavantaje. Nu există garanție că nu s-a strecurat pe undeva vreo greșală. Trebuie să știi matematică dacă vrei să îl folosești în siguranță. Nu este prietenos, așa cum este de exemplu Excel-ul. Ca să poți folosi help-ul există chiar o instrucțiune `help.search`; nu are un help ușor accesibil. Dar nici matematica nu este ușor accesibilă, iar cartea aceasta este destinată unor cititori cu pregătire matematică – fie ei studenți la matematică, informatică, politehnică sau construcții ⁵.

Cartea este intenționată a fi o inițiere în R. Sistemul R este imens – există deja mai mult de 1000 de pachete. Noi vrem să prezentăm doar pachetul de bază, acela care se descarcă de la CRAN (<http://cran.r-project.org>). Acesta este suficient pentru a putea rula R și conține funcțiile fundamentale ale limbajului.

Utilizatorul poate fi descurajat de cantitatea imensă de documentație care există în legătură cu R. Numai editura Springer a publicat, începând cu 2004, câteva zeci de cărți dedicate acestui limbaj, scrise pentru cititori cu pregătire matematică. O listă aproape completă se poate găsi la <http://www.r-project.org/doc/bib/R-books.html> ⁶.

Cunoscătorii apreciază R-ul pentru facilitățile sale grafice ⁷.

Pentru un începător în R există cel puțin trei mari dificultăți.

Prima este legată de sintaxă. O paranteză greșit pusă, sau o paranteză `[.]` în loc de `(.)` sau de `{.}` poate face necazuri foarte mari. Mesajele de eroare sunt sibilnice.

A doua este legată de scripturi. După cum se va vedea, utilizatorul poate da instrucțiunile ori în linie de comandă, ori în pachete de instrucțiuni numite scripturi. Un script este un fișier cu extensia `.R`. Ca să deschizi un script, însă trebuie să știi unde e salvat. Un script nu este un program sau o singură subrutină, el poate să conțină mai multe asemenea programe. Pentru a-l rula, este suficient să îl deschizi, să cauți paragraful care te interesează, să îl selectezi și apoi să tastezi `Ctrl+R`. Se va executa exact ceea ce dorești, dar asta se învață mai greu. Pe de altă parte, este bine ca proiectele diferite să fie puse în directoare diferite, fiindcă R păstrează în memorie variabilele ultimei sesiuni de lucru și probabilitatea de a da același nume unor variabile diferite este foarte mare.

În fine, a treia dificultate este legată de introducerea datelor din fișiere. Începătorul poate fi foarte frustrat încercând să copieze un tabel din Excel într-o matrice din R: dacă tabelul din Excel nu este o matrice – de exemplu dacă are celule lăsate necompletate atunci R dă mesaje de eroare neinteligibile cum ar fi „Error in scan (file, what, nmax, sep, dec, quote, skip, nlines, na.strings, : line 3 did not have 3 elements”. Același lucru îl putem spune și despre salvarea în fișiere a rezultatelor obținute prin calcule: a transfera un vector din R în Excel este o operație care poate provoca mare frustrare utilizatorului neavertizat, așa cum se va vedea la capitolul doi.

⁵ Chiar și persoane fără o pregătire matematică specială pot folosi sistemul R dacă au un minimum de pregătire statistică, de exemplu studenți la biologie, sociologie, agronomie, etc.

⁶ Manualul complet “R Reference Manual” care este conținut în kitul de instalare R are 1884 de pagini din care 60 pagini reprezintă indexul cu peste 4000 de funcții puse la dispoziția utilizatorilor care știu să le folosească.

⁷ Murrell, P. (2005) *R Graphics*. Chapman & Hall/CRC Press.

Odată depășite aceste dificultăți inerente începutului, apar și satisfacțiile. De exemplu a găsi minimul, maximul, mediana unui șir de date, ai face histograma, a-l sorta devin o joacă de copii.

În primele patru capitole se prezintă principalele posibilități de utilizare a mediului R în regim de comandă. În capitolul al cincilea se prezintă posibilitățile de programare în R, iar în capitolul al șaselea sunt prezentate aplicații programate în R din diverse domenii care au legătură cu statistica. În ultimul capitol se prezintă teste statistice și serii de timp.

Maria Miroiu a scris capitolele 1,2 și 7.1; Viorel Petrehuș a scris capitolele 3,4,5 și 7.2 iar Gheorghiță Zbăganu a scris capitolul 6.

Variantele de R utilizate au fost: 2.12 și 2.14.

Accentuăm: cartea se adresează utilizatorilor începători ai sistemului R, dar care au o minimă pregătire în matematică, îndeosebi în domeniul statisticii. Se presupune că există la cititor o oarecare familiaritate cu termenii de variabilă aleatoare, repartiție, test statistic. Este important să accentuăm însă că sistemul R poate fi utilizat în mai multe domenii ale matematicii aplicate, nu numai în statistică, de exemplu pentru rezolvarea ecuațiilor și sistemelor de ecuații diferențiale.

În final avem de adăugat următoarele: cartea a fost elaborată în cadrul proiectului POSDRU/56/1.2/S/32768, “Formarea cadrelor didactice universitare și a studenților în domeniul utilizării unor instrumente moderne de predare-învățare-evaluare pentru disciplinele matematice, în vederea creării de competențe performante și practice pentru piața muncii”.

Finanțat din Fondul Social European și implementat de către Ministerul Educației, Cercetării, Tineretului și Sportului, în colaborare cu The Red Point, Oameni și Companii, Universitatea din București, Universitatea Tehnică de Construcții din București, Universitatea „Politehnica” din București, Universitatea din Pitești, Universitatea Tehnică „Gheorghe Asachi” din Iași, Universitatea de Vest din Timișoara, Universitatea „Dunărea de Jos” din Galați, Universitatea Tehnică din Cluj-Napoca, Universitatea “1 Decembrie 1918” din Alba-Iulia, proiectul contribuie în mod direct la realizarea obiectivului general al Programului Operațional Sectorial de Dezvoltare a Resurselor Umane – POSDRU și se înscrie în domeniul major de intervenție 1.2 Calitate în învățământul superior.

Proiectul are ca obiectiv adaptarea programelor de studii ale disciplinelor matematice la cerințele pieței muncii și crearea de mecanisme și instrumente de extindere a oportunităților de învățare.

Obiectivele de interes ale proiectului sunt evaluarea nevoilor educaționale ale cadrelor didactice și studenților legate de utilizarea matematicii în învățământul superior, masterate și doctorate precum și analizarea eficacității și relevanței curriculelor actuale la nivel de performanță și eficiență, în vederea dezvoltării de cunoștințe și competențe pentru studenții care învață discipline matematice în universități. Obiectivele specifice la care răspunde materialul de față sunt: dezvoltarea și armonizarea curriculelor disciplinelor matematice, conform exigențelor moderne; elaborarea și implementarea unui program de formare a cadrelor didactice și a studenților interesați din universitățile partenere, bazat pe dezvoltarea și armonizarea de curriculum; crearea unei baze de resurse inovative și funcționale pentru predarea și învățarea în disciplinele matematice pentru învățământul universitar. sunt obiectivele specifice care au ca raspuns materialul de față.

Formarea de competențe cheie de matematică și informatică presupune crearea de abilități de care fiecare individ are nevoie pentru dezvoltarea personală, incluziune socială și inserție pe piața muncii. Se poate constata însă că programele disciplinelor de matematică nu au întotdeauna în vedere identificarea și sprijinirea elevilor și studenților potențial talentați la matematică.

În acest context, analiza flexibilității curriculei, însoțită de analiza metodelor și instrumentelor folosite pentru identificarea și motivarea studenților talentați la matematică ar putea răspunde deopotrivă cerințelor de masă, cât și celor de elită.

Viziunea pe termen lung a acestui proiect preconizează determinarea unor schimbări în abordarea fenomenului matematic pe mai multe planuri: informarea unui număr cât mai mare de membri ai societății în legătură cu rolul și locul matematicii în educația de bază în instrucție și în descoperirile științifice menite să îmbunătățească calitatea vieții, inclusiv popularizarea unor mari descoperiri tehnice, și nu numai, în care matematica cea mai avansată a jucat un rol hotărâtor. De asemenea, se urmărește evidențierea a noi motivații solide pentru învățarea și studiul matematicii la nivelele de bază și la nivel de performanță; stimularea creativității și formarea la viitorii cercetători matematicieni a unei atitudini deschise față de însușirea aspectelor specifice din alte științe, în scopul participării cu succes în echipe mixte de cercetare sau a abordării unei cercetări inter și multi disciplinare; identificarea unor forme de pregătire adecvată de matematică pentru viitorii studenți ai disciplinelor matematice, în scopul utilizării la nivel de performanță a aparatului matematic în construirea unei cariere profesionale.

CUPRINS

Introducere	1
Cap.1 Preliminarii	5
Cap.2 Tipuri de date	10
Cap.3 Funcții predefinite în R	30
Cap.4 Funcții grafice în R	52
Cap.5 Programare în R	64
Cap.6 Aplicații diverse programate în R	77
Cap.7 Aplicații rezolvate prin funcții predefinite în R	
7.1 Teste statistice	122
7.2 Serii de timp	151
Bibliografie	179

Capitolul 1

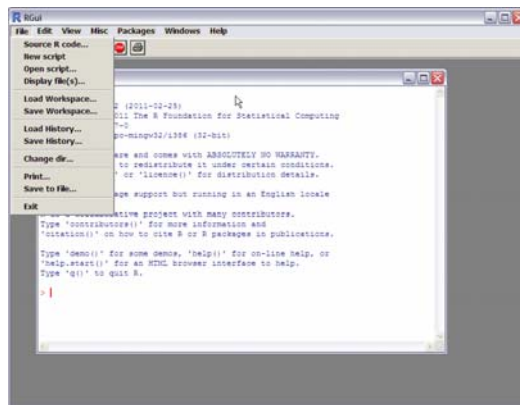
Preliminarii

R este un mediu integrat de facilități software destinat în special analizei statistice și reprezentării grafice a datelor. Acesta include un limbaj simplu de programare, limbajul S, care permite utilizarea pachetelor existente și construirea de pachete noi de programe destinate analizei statistice și reprezentării grafice a datelor.

R este un mediu de dezvoltare gratuit distribuit sub licență publică generală GNU. A fost inițial scris de Ross Ihaka și Robert Gentleman de la departmentul de statistică de la Universitatea din Auckland, New Zealand. Din 1997 s-a format un grup principal ("R Core Team") care poate modifica variante ale R-ului. R-ul s-a dezvoltat rapid și s-a extins la o mare colecție de pachete de programe.

Site-ul oficial de unde se pot obține gratuit versiuni compilate de R pentru OS X, Linux și Windows este <http://www.r-project.org>. Tot aici se pot găsi și resurse pentru instalarea și utilizarea R-ului. Cea mai recentă versiune este 2.13.2 disponibilă din 30.09.2011. Odată cu R se descarcă și o serie de pachete de programe (numite și pachete standard), dar sunt disponibile și alte pachete de programe R prin intermediul CRAN (Comprehensive R Archive Network – <http://CRAN.R-project.org>) care conține diferite variante ale sistemului R și ale pachetelor de programe R.

Kitul de instalare R vine cu manuale în care este descris în detaliu modul de funcționare al programului. De asemenea se găsesc cărți specializate pe diverse domenii unde se utilizează statistica și unde se arată cum se poate utiliza R pentru calculele necesare.



La lansarea în execuție a R-ului apare o fereastră ca cea alăturată și se poate lucra direct cu R în fereastra consolă interioară.

În meniul „Packages” se găsesc opțiuni de instalare separată a unor pachete suplimentare de programe specializate pe diverse domenii prin care posibilitățile lui R se pot extinde foarte mult.

Lucrul cu R se poate face

- fie în regim interactiv: în fereastra Rconsole, după prompter-ul `>`, se scrie o comandă acceptată de R, apoi se apasă tasta Enter și rezultatul apare imediat mai jos în aceeași fereastră,
- fie prin script-uri, în regim de programare: comenzile de calcul sunt grupate într-un fișier cu extensia `.r` care poate fi apelat în mai multe moduri. Un fișier script poate fi scris în orice editor de text, dar R are un editor ce se apelează din meniul „File” prin „Open script...” sau „New script”. Fișierul script poate fi executat după ce se încarcă în mediul R prin *File* → *Open script...* sau se poate folosi direct funcția `source()`, de exemplu, presupunând că fișierul script se numește `test.r` și se află în directorul de lucru curent, comanda scrisă la prompt: `>source("test.r")` va avea ca efect încărcarea în spațiul de lucru a fișierului și executarea comenzilor conținute în fișierul „test.r”.

R dispune de un mecanism de rechemare a comenzilor executate în linia de comandă în fereastra RConsole. Prin apăsarea tastei de deplasare pe verticală în sus (↑), se pot obține comenzile scrise anterior. Acestea pot fi reexecutate sau modificate și apoi executate.

Toate datele citite sau salvate prin comenzi R sunt asociate cu directorul curent care poate fi schimbat prin *File* → *Change dir...* la calea dorită de utilizator.

Directorul curent (de lucru) poate fi aflat prin comanda `getwd()` și poate fi fixat și prin comanda `setwd("directorul de lucru dorit")`.

R dispune de funcții care permit obținerea rapidă de informații referitoare la funcțiile și facilitățile sistemului R instalat pe calculatorul propriu. Pentru a obține informații despre facilitățile oferite de funcția `help()`, din R se execută comanda

```
> help.start()
```

Se va lansa o pagină web cu o serie link-uri care permit accesarea informațiilor legate de sistemul R.

Pentru a obține informații despre o anumită funcție R, de exemplu `plot`, se execută comanda

```
> help(plot)
```

sau comanda alternativă

```
> ?plot
```

Aceasta va deschide un Web browser cu o pagină în care se afișează informații referitoare la funcția `plot`.

Pentru funcțiile specificate prin cuvinte cheie, ca de exemplu, `if`, `while`, `function`, argumentele trebuie incluse între ghilimele:

```
> help("while")
```

Pentru a afișa fișiere, titluri, obiecte, etc. ce conțin un anumit șir de caractere, se folosește comanda

```
> help.search(șir de caractere, ...)
```

De exemplu:

```
> help.search("poisson")
```

sau echivalent

```
> ??poisson
```

conduce la afișarea următoarelor rezultate

```
Help files with alias or concept or title matching  
'poisson' using fuzzy matching:
```

<code>boot::poisons</code>	Animal Survival Times
<code>stats::family</code>	Family Objects for Models
<code>stats::Poisson</code>	The Poisson Distribution
<code>stats::poisson.test</code>	Exact Poisson tests

Dacă dorim să încărcăm un anumit pachet de programe R, mai întâi trebuie descărcat în subdirectorul library al directorului R (de exemplu în: *C:\Program Files\R\R-2.13.0\library*). Pentru încărcarea de exemplu a pachetului *graphics* (the R graphics package), putem folosi funcția

```
> library(graphics)
```

Dacă dorim să aflăm informații despre un anumit pachet de programe, de exemplu utils (the R Utils package)

```
> library(help="utils").
```

Pentru a vedea ce pachete de programe R sunt disponibile în versiunea curentă R, se poate folosi funcția

```
> library()
```

De exemplu, pachetele disponibile în versiunea R-2.13.0 sunt:

base	The R Base Package
boot	Bootstrap R (S-Plus) Functions (Canty)
class	Functions for Classification
cluster	Cluster Analysis Extended Rousseeuw et al.
codetools	Code Analysis Tools for R
compiler	The R Compiler Package
datasets	The R Datasets Package
foreign	Read Data Stored by Minitab, S, SAS, SPSS, Stata, Systat, dBase, ...
graphics	The R Graphics Package
grDevices	The R Graphics Devices and Support for Colours and Fonts
grid	The Grid Graphics Package
KernSmooth	Functions for kernel smoothing for Wand & Jones (1995)
lattice	Lattice Graphics
MASS	Support Functions and Datasets for Venables and Ripley's MASS
Matrix	Sparse and Dense Matrix Classes and Methods
methods	Formal Methods and Classes
mgcv	GAMs with GCV/AIC/REML smoothness estimation and GAMMs by PQL
nlme	Linear and Nonlinear Mixed Effects Models
nnet	Feed-forward Neural Networks and Multinomial Log-Linear Models
rpart	Recursive Partitioning
spatial	Functions for Kriging and Point Pattern Analysis
splines	Regression Spline Functions and Classes
stats	The R Stats Package
stats4	Statistical Functions using S4 Classes
survival	Survival analysis, including penalised likelihood.
Tcltk	Tcl/Tk Interface
tools	Tools for Package Development
utils	The R Utils Package

Pentru a vizualiza exemple referitoare la anumite funcții R, de exemplu pentru funcția `mean`, se poate proceda astfel:

```
> example(mean)
```

iar răspunsul sistemului este următorul:

```
mean> x <- c(0:10, 50)
mean> xm <- mean(x)
mean> c(xm, mean(x, trim = 0.10))
[1] 8.75 5.50
mean> mean(USArrests, trim = 0.2)
Murder  Assault UrbanPop  Rape
7.42    167.60    66.20    20.16
```

Odată cu lansarea R-ul se deschide o nouă sesiune R. Ieșirea din sesiunea R se poate face utilizând comanda:

```
> q()
```

sau echivalent

```
> quit()
```

Pe ecran va apare o întrebare referitoare la memorarea sesiunii de lucru (instrucțiuni lansate la linia de comandă, obiecte cu care s-a operat, etc.) care se poate salva într-un fișier și redeschide mai târziu sau se poate renunța definitiv la aceasta.

Dacă se dorește salvarea spațiului de lucru (ce cuprinde toate variabilele create în R), într-un fișier anume se poate folosi comanda `save`. Să presupunem de exemplu că într-o nouă sesiune R, în spațiul de lucru, printre altele, avem și variabila `x`:

```
> x=2          #se creeaza si initializeaza varibila x
> save(x, file="C:/Lucru/date.RData")
> # se salveaza variabila x in fisierul date.Rdata
> q()          #părăsirea sesiunii R
```

Funcția `save` este echivalentă cu **`saveimage`**.

Să mai precizăm că R este **case sensitive** și deci `x` și `X` sunt interpretate ca simboluri diferite și se vor referi la variabile diferite.

De asemenea, să precizăm că orice șir de caractere precedat de simbolul '#' este interpretat ca fiind un **comentariu**.

La redeschiderea unei noi sesiuni R se poate proceda în modul următor:

```
> load("c:/lucru/date.RData")
> #se incarca obiectele din fisierul date
> x          #se afiseaza valoarea variabilei x
[1] 2        #raspunsul sistemului
```

În mod similar, executând dublu click pe fișierul „date.Rdate”, se realizează deschiderea unei noi sesiuni R și apoi încărcarea fișierului „date.Rdate”. De asemenea, salvarea spațiului de lucru curent se poate realiza prin meniul *File* → *Save Workspace* ..., respectiv încărcarea unui fișier Rdate se poate realiza prin meniul *File* → *Load Workspace* ...

Comenzile se pot separa în R fie prin simbolul punct și virgulă fie se pot scrie pe linii diferite. De exemplu:

```
> x1=1; x2=3.7;
> x3=5.1;
> x1+x2
[1] 4.7
> x1-x3
[1] -4.1
```

Dacă o comandă nu este completă la sfârșitul unei linii de la prompt-ul >, R va lansa un nou prompt, de obicei simbolul '+' în linia sau liniile următoare. Acest simbol dispare când comanda este completă. De exemplu:

```
> x=4
> y=6
> z=x*
+ y
> z
[1] 24
```

Istoria comenzilor introduse la prompt în Rconsole într-o sesiune de lucru R poate fi salvată, iar apoi reîncărcată folosind comenzile **savehistory** și **loadhistory**, care se regăsesc și în meniu: *File* → *Save History...*, respectiv *File* → *Load History...*. Fișierele în care se salvează / de unde se încarcă comenzile au extensia .Rhistory.

Pentru alte comenzi / funcții R de bază se poate consulta help-ul Pachetul R Utils:

```
> library(help="utils")
```

Capitolul 2

Tipuri de date

Entitățile cu care operează R pentru a memora și a lucra cu date se numesc **tipuri de date** sau **obiecte**. Acestea pot fi: scalari, vectori, matrice, liste, tablouri, șiruri de caractere, factori, dataframe-urile, etc. De asemenea, se pot defini tipuri noi de obiecte.

Funcția **typeof**(obiect) sau echivalent **mode**(obiect) pot fi folosite pentru a afla tipul atributelor intrinseci ale unui obiect. Câteva dintre cele mai des folosite tipuri de obiecte: "numeric" (întreg="integer" sau real="double"), "complex", "logical" (logic), "character" (caracter) sau "raw" (vector de biți, fiecare componentă a acestuia fiind reprezentată separat ca pereche de cifre hexazecimale). De exemplu:

```
> x=3.2          > s="unu"          > s="sir"
> typeof(x)      > typeof(x)      > b=charToRaw(s);b
[1] "numeric"     [1] "character"    [1] 73 69 72
                  > typeof(b)
                  [1] "raw"
```

```
> x=1+2i          > e=1<3
> typeof(x)      > typeof(e)
[1] "complex"     [1] "logical"
```

Accesul la obiecte se face prin intermediul unor simboluri numite **variabile**. În R, variabilele sunt ele însele obiecte și pot fi manipulate în același mod ca și obiectele. Numele variabilele din R încep cu o literă, apoi pot conține litere, cifre, „.” sau „_”.

Pentru a atribui o valoare unei variabile se poate folosi **operatorul de atribuire** („=” sau „<=” sau „<<=” sau „->” sau „->>”) sau se poate folosi funcția **assign** din pachetul de programe R „base”:

```
> variabila = valoare
> variabila <- valoare
> variabila <<- valoare
> valoare -> variabila
> valoare ->> variabila
> assign("variabila",valoare)
```

Se observă faptul că după atribuire, rezultatul nu este afișat automat.

Afișarea valorii unei variabile se face explicit, prin apelarea numelui variabilei sau prin folosirea funcției **print**. De exemplu:

```
> x=1.5
> x
[1] 1.5
> print(x)
[1] 1.5
```

La afișarea valorilor variabilelor se poate folosi funcția **options**(nume=valoare) pentru a stabili diverse opțiuni. De exemplu, pentru a stabili numărul de zecimale la 2 se apelează înainte de afișare: **options(digits=2)**. Valorile valide sunt 1..22 cu valoarea implicită 7. Însă acest număr de zecimale este doar sugestie. De exemplu:

```
> options(digits=5)
> x=1.3
```

```
> x
[1] 1.3
> y=1/3; y
[1] 0.33333
```

2.1. Numere

R lucrează cu numere întregi, reale sau complexe. Numerele reale se scriu cu punct zecimal sau cu exponent: $0.001 = 1.e-3$.

Pentru a calcula o expresie numerică, se poate edita în RConsole de la prompter expresia corespunzătoare, apoi se tastează <Enter>. De exemplu,

```
>> 2*6
[1] 12          #rezultat
```

Cea de-a doua linie reprezintă răspunsul sistemului R, iar [1] indică numărul de ordine al primei valori afișate pe acest rând. Alt exemplu:

```
>> sqrt(2)
[1] 1.414214
```

Pentru a genera 14 numere aleatoare distribuite uniform în intervalul (0,1) se poate folosi comanda:

```
> runif(n=14, min=0, max=1)
[1] 0.70871023 0.39347162 0.51954422 0.23446865 0.38287326 0.26328696
[7] 0.56719362 0.82779636 0.01660613 0.49667370 0.20340420 0.96901185
[13] 0.25215480 0.94807566
```

Operatorii binari cu care poate lucra R sunt:

Operația matematică	Simbolul R	Expresia R
Adunarea $x+y$	+	$x+y$
Scăderea $x-y$	-	$x-y$
Înmulțirea xy	*	$x*y$
Împărțirea x/y	/	x/y
Exponențiala x^y	^	x^y
Câtul întreg al împărțirii lui x la y : $[x/y]$	%/%	$x\%/%y$
Restul împărțirii lui x la y : $x - y * [x/y]$	%%	$x\%y$

Expresia matematică se evaluează de la stânga la dreapta, ținând cont de și ordinul de prioritate:

Ordinul de prioritate	
1	Paranteza
2	Exponențiala
3	Înmulțirea, împărțirea, câtul, restul
4	Adunarea și scăderea

Operatorii %% și %/% pot fi folosiți și pentru operatori neîntregi. De exemplu:

```
> 10.2%/%1.2
[1] 8
> 10.2%%1.2
[1] 0.6
```

În R există constante speciale precum: **Inf** (ce ține locul lui infinit ∞), **NaN** (“not-a-number”, ce indică un rezultat numeric nedefinit ca $0/0$, $(\pm\infty)/(\pm\infty)$, $\infty - \infty$ sau $0 \times (\pm\infty)$), **NULL** (indică un obiect gol) și **NA** (“not available”). De exemplu:

```
> 1/0
[1] Inf
> 5^987
[1] Inf
> 0/0
[1] NaN
> Inf/Inf
[1] NaN
> u=c(1,2,NA,4);u #valoare lipsa reprezentata prin NA
[1] 1 2 NA 4
> A = matrix(data=NA,nrow=2,ncol=3); A
      [,1] [,2] [,3]
[1,]  NA  NA  NA
[2,]  NA  NA  NA
```

2.2. Vectori

Vectorii sunt structuri de date cu una sau mai multe valori de același tip: întreg, real, complex, logic, caracter sau raw (șiruri de biți). Vectorii cu o singură componentă se numesc *scalari*.

În R, nu sunt permisi indici negativi, iar indicele 0 este ignorat. Așadar, indicii acceptați sunt 1, 2, ...

Vectorii pot fi creați interactiv sau se pot importa din fișiere. Crearea interactivă a vectorilor se poate face utilizând funcția de concatenare „c”. De exemplu:

```
> u = c(1,2,3)
> p = c(1.25, 0.75)
> v = c(1+1i, 1-2i)
```

Pentru vizualizarea unei componente, de exemplu a doua, se poate face astfel:

```
> u[2]
[1] 2
```

Funcția de concatenare se poate folosi și pentru a construi vectori de dimensiune mai mare prin concatenarea mai multor vectori de dimensiuni mai mici. În exemplul următor s-a creat un vector numeric, x_1 , un vector x_2 cu elementele de tip șir de caractere și s-au combinat cei doi vectori în unul singur, x_3 . Vectorul x_3 trebuie să aibă componentele de același tip și tipul este șir de caractere. *Regula de conversie a datelor este logice => numere => caractere.*

```
> x1=c(1,2,3)           #vector numeric
> x2=c("unu","doi")     #vector de siruri
> x3=c(x1,x2); x3        #concatenarea celor doi vectori
[1] "1" "2" "3" "unu" "doi"  #conversie la siruri
```

Limbajul R permite efectuarea interactivă a diferitor operații cu vectori numerici, precum adunarea (+), scăderea (-), înmulțirea(*), împărțirea (/) și ridicarea la putere (^). De exemplu, pentru vectorii $x = (2\ 4\ 6)$ și $y = (1\ 2\ 3)$ avem:

```
> x=c(2,4,6)      > x*y          > x+2
> y=c(1,2,3)      [1] 2 8 18       [1] 4 6 8
> x+y             > x/y          > options(digits=2)
[1] 3 6 9          [1] 2 2 2
> x-y             > x^y          > 1/x
[1] 1 2 3          [1] 2 16 216      [1] 0.50 0.25 0.17
```

Un operator special este `'%*%'` care aplicat unor vectori conduce la determinarea produsului scalar, ca în exemplul de mai jos:

```
> x=c(2,4,6)
> y=c(1,2,3)
> x%*%y           #produsul scalar a vectorilor x si y
               [,1]
[1,]          28
```

R dispune de o serie de facilități pentru generarea secvențelor de numere, organizate ca numere sau matrice. De exemplu, pentru a genera secvența de numere consecutive de la 1 la 5 și memorarea acesteia în vectorul u se poate proceda astfel:

```
> u=1:5; u
[1] 1 2 3 4 5
```

În mod similar, pentru se poate genera secvența de numere consecutive de la 5 la 1 și apoi memorarea acesteia în vectorul v folosind instrucțiunea

```
> v=5:1; v
[1] 5 4 3 2 1
```

Generarea de numere se poate face și cu funcția `seq(from,to,by)`. Primele două argumente specifică valorile de început (from) și sfârșit (to) ale secvenței, iar ultimul argument (by) reprezintă pasul. De exemplu:

```
> v1=seq(3,5)      #sau v1=seq(from=3,to=5)
[1] 3 4 5
> v2=seq(from=1, to=2, by=0.2); v2
[1] 1.0 1.2 1.4 1.6 1.8 2.0
> seq(4,16,length=6)
[1] 4.0 5.2 6.4 7.6 8.8 10.0
```

O funcție înrudită cu `seq()` este `rep()`. Aceasta se utilizează pentru a replica un obiect în diferite moduri.

```
> x=c(1,2,3)
> x1=rep(x, times=2); x1      #x se repeta de 2 ori
[1] 1 2 3 1 2 3
> x2=rep(x, each=2); x2
#fiecare componenta a lui x se va repeta de 2 ori
[1] 1 1 2 2 3 3
> rep(1:3,c(2,4,6))
[1] 1 1 2 2 2 2 3 3 3 3 3 3
```

#prima componenta s-a multiplicat de 2 ori, a doua de 4 ori etc

```
> rep(0,10)
[1] 0 0 0 0 0 0 0 0 0 0
```

Este important să putem identifica anumite șiruri de numere. Operatorul „%in%” care testează existența unei secvențe într-un vector de numere este deosebit de util în acest scop. De exemplu:

```
> x=rep(1:3,rep(2,3)); x
[1] 1 1 2 2 3 3
> x[x %in% c(1,3)]
#se testează existența valorilor 1 și 3
[1] 1 1 3 3
```

Pentru a determina lungimea unui vector se poate folosi funcția length, ca mai jos:

```
> x=seq(-15,15,by=5); x
[1] -15 -10 -5 0 5 10 15
> length(x)          #lungimea vectorului x
[1] 7
```

Elementele unui **vector logic** sunt succesiuni de TRUE (adevărat), FALSE (fals) și NA („not available”). Deseori, TRUE se abreviază cu T, iar FALSE cu F, dar este posibil ca variantele abreviate să fie variabile deja definite.

Vectorii logici sunt generați atunci când se impun anumite condiții. De exemplu, dacă punem condiția ca elementele vectorului $v = (2 \ -1 \ 5)$ să fie pozitive, putem genera vectorul logic (TRUE, FALSE, TRUE) cu secvența de instrucțiuni:

```
> v=c(2,-1,5); v
[1] 2 -1 5
> y=v>0
> y
[1] TRUE FALSE TRUE
```

Primul și ultimul element al lui v sunt mai mari decât zero și condiția este îndeplinită (TRUE), iar al doilea este negativ și condiția nu este îndeplinită (FALSE).

Operatorii relaționali sunt: < (‘mai mic’), <= (‘mai mic sau egal’), > (‘mai mare’), >= (‘mai mare sau egal’), == (‘egalitate’) și != (diferit).

Operatorii logici: & (și), | (sau), ! (negarea).

Vectorii logici pot fi utilizați și în operații aritmetice. În acest caz, ei se transformă în vectori numerici, FALSE devenind 0, iar TRUE devenind 1:

```
> p=c(TRUE,FALSE,TRUE)
> q=c(TRUE,FALSE,FALSE)
> p
[1] TRUE FALSE TRUE
> q
[1] TRUE FALSE FALSE
> p+q
[1] 2 0 1
```

Vectorii de tip caracter sunt secvențe de caractere delimitate de ghilimele. Aceștia pot fi concatenați utilizând funcția `c()`. De asemenea, vectorii de tip caracter pot fi concatenați folosind funcția `paste()`, care acceptă un număr arbitrar de argumente pe care le concatenează unul câte unul, separându-le printr-un caracter blank. Orice număr este transformat în caracter.


```

> nume="Popescu"
> prenume="Ion"
> numeintreg=paste(nume,prenume)
> numeintreg
[1] "Popescu Ion"
> n=c(nume,prenume)
> n
[1] "Popescu" "Ion"
> n[1]
[1] "Popescu"
> n[2]
[1] "Ion"

```

Funcția `strsplit` se poate folosi pentru extragerea unor caractere dintr-un vector utilizat ca argument. De exemplu, pentru eliminarea caracterului 'a' din șirul "Transilvania" se poate proceda astfel

```

> strsplit("Transilvania","a")
[[1]]
[1] "Tr" "nsilv" "ni"

```

Funcția `sort` se poate folosi pentru a rearanja / sorta elementele unui vector caracter în ordine alfabetică. De exemplu:

```

> s=c("Iulia","Magda","Ana","Paula")
> s
[1] "Iulia" "Magda" "Ana"   "Paula"
> sort(s)
[1] "Ana"   "Iulia" "Magda" "Paula"

```

2.3. Matrice și tablouri de numere

Matricele pot fi **create** în diverse moduri. Modalitatea cea mai simplă este de a utiliza vectori de aceeași lungime care se pun ca și coloane ale matricei prin instrucțiunea `cbind()` sau ca și linii ale matricei prin instrucțiunea `rbind()`, ca mai jos. **Dimensiunile** unei matrice se pot determina prin folosirea funcției `dim`.

```

> y1=c(1,2,3)
> y2=c(4,5,6)
> A=cbind(y1,y2) # concatenare pe orizontala
> B=rbind(y1,y2) # concatenare pe verticala
> A
      y1 y2
[1,]  1  4
[2,]  2  5
[3,]  3  6
> B
      [,1] [,2] [,3]
y1      1    2    3
y2      4    5    6
> dim(B)
[1] 2 3

```

Pentru ca doi vectori să se poată concatena pe orizontală aceștia trebuie să aibă același număr de linii, iar pentru se putea concatena pe verticală aceștia trebuie să aibă același număr de coloane.

Funcțiile `cbind` și `rbind` se pot aplica și la matrice. De exemplu:

```
> a<-c(-1,0,1)
> b<-c(-2,-1,1)
> c<-cbind(a,b); #concatenare vectori pe orizontala
> c
      a  b
[1,] -1 -2
[2,]  0 -1
[3,]  1  1

> d<-cbind(c,c);d #concatenare matrice pe orizontala
      a  b  a  b
[1,] -1 -2 -1 -2
[2,]  0 -1  0 -1
[3,]  1  1  1  1

> d=rbind(d,d); d #concatenare matrice pe verticala
      a  b  a  b
[1,] -1 -2 -1 -2
[2,]  0 -1  0 -1
[3,]  1  1  1  1
[4,] -1 -2 -1 -2
[5,]  0 -1  0 -1
[6,]  1  1  1  1

> e=cbind(5,c) #se adauga o coloana de 5 la matricea c
> e
      a  b
[1,]  5 -1 -2
[2,]  5  0 -1
[3,]  5  1  1

> f=rbind(7,c);f #se adauga o linie de 7 la matricea c
      a  b
[1,]  7  7
[2,] -1 -2
[3,]  0 -1
[4,]  1  1
```

Dacă tipurile vectorilor care intră în matrice sunt diferite atunci se aplică regula de conversie `logice=>numere=>caractere`. Este posibil ca lungimile vectorilor să fie diferite, dar cea mai mare lungime să fie multiplu întreg al lungimilor mai mici, caz în care vectorii de lungime mai mică se repetă.

O a doua modalitate de a crea matrice, cât și tablouri multidimensionale (arrays) se poate face și prin comenzile `matrix` și `array` plecând de la un vector prin aranjarea elementelor lui într-un tablou de dimensiuni date.

```
> y=1:12
> m=matrix(y, nrow=4, ncol=3) #completare pe coloane
> m
      [,1] [,2] [,3]
[1,]    1    5    9
[2,]    2    6   10
[3,]    3    7   11
[4,]    4    8   12
```

Se observă că din elementele vectorului y_3 se formează coloanele matricii m_3 . Numărul de elemente ale primului parametru trebuie să fie multiplu sau submultiplu al lui `nrow*ncol`. În cazul în care este submultiplu, valorile vor fi ciclate până când toate valorile sunt completate. De exemplu:

```
> m=matrix(1:6, nrow=4, ncol=3); m
> #completare pe coloane
      [,1] [,2] [,3]
[1,]     1     5     3
[2,]     2     6     4
[3,]     3     1     5
[4,]     4     2     6
```

sau

```
> m=matrix(1:14, nrow=4, ncol=3)
Warning message:
In matrix(1:14, nrow = 4, ncol = 3) :
data length [14] is not a sub-multiple or multiple
of the number of rows [4]
```

Dacă dorim ca din elementele lui să formăm liniile, mai adăugăm indicația `byrow=TRUE` (valoarea implicită fiind `byrow=FALSE`). Astfel se obține

```
> m4=matrix(1:12, nrow=4, ncol=3, byrow=TRUE)
> #completare pe linii
> m4
      [,1] [,2] [,3]
[1,]     1     2     3
[2,]     4     5     6
[3,]     7     8     9
[4,]    10    11    12
```

O altă modalitate de a crea matrice este cea care convertește efectiv un vector în matrice, folosind funcția `dim` ca în exemplul de mai jos:

```
> x=seq(1:12); x      #x este vector
[1]  1  2  3  4  5  6  7  8  9 10 11 12

> dim(x)=c(4,3); x    #x este matrice
      [,1] [,2] [,3]
[1,]     1     5     9
[2,]     2     6    10
[3,]     3     7    11
[4,]     4     8    12
```

Accesul la elementele din vector, matrice sau array se face punând indicii de poziție ai elementului între paranteze drepte, și eventual separându-i prin virgulă:

```
> m=matrix(1:6, nrow=2, ncol=3)
> m      # sau m[, ]
      [,1] [,2] [,3]
[1,]     1     3     5
[2,]     2     4     6

> m[1,3]
[1] 5
> m[2,] #sau m[2,1:3] (linia a doua)
[1] 2 4 6
> m[2,1:2] #linia a doua, primele doua elemente
[1] 2 4
```

```

> m[,1] #sau m[1:2,1] (coloana 1)
[1] 1 2

> m[,2:3] #coloanele 2 si 3
      [,1] [,2]
[1,]     3     5
[2,]     4     6

```

Pentru a prezenta **operațiile cu matrice**, să considerăm întâi matricele

$$A = \begin{pmatrix} 2 & 4 \\ 1 & 10 \\ 6 & 0 \end{pmatrix} \text{ și } B = \begin{pmatrix} 4 & 6 \\ 8 & 5 \\ 1 & 2 \end{pmatrix} :$$

```

> A=matrix(c(2,4,1,10,6,0),ncol=2,byrow=TRUE)
> A
      [,1] [,2]
[1,]     2     4
[2,]     1    10
[3,]     6     0

> B=matrix(c(4,6,8,5,1,2),ncol=2,byrow=TRUE)
> B
      [,1] [,2]
[1,]     4     6
[2,]     8     5
[3,]     1     2

```

Apoi

```

> A+B
      [,1] [,2]
[1,]     6    10
[2,]     9    15
[3,]     7     2

> A-B
      [,1] [,2]
[1,]    -2    -2
[2,]    -7     5
[3,]     5    -2

> A*B
> # (aij*bij)i,j = inmultirea element cu element
      [,1] [,2]
[1,]     8    24
[2,]     8    50
[3,]     6     0

> A/B      # (aij/bij)i,j
> # (aij/bij)i,j = impartirea element cu element
      [,1] [,2]
[1,] 0.500 0.6666667
[2,] 0.125 2.0000000
[3,] 6.000 0.0000000

> A^B      # (aij^bij)i,j
> # (aij^bij)i,j=ridicarea la putere element cu elem.
      [,1] [,2]
[1,]    16   4096
[2,]     1 100000
[3,]     6     0

```

De observat că înmulțirea matricelor folosind operatorul '*' nu este înmulțirea clasică ci înmulțirea element cu element. Analog, folosirea operatorului '/', respectiv '^' conduce la construirea unei matrice în care fiecare element este de tipul a_{ij} / b_{ij} , respectiv $a_{ij} ^ b_{ij}$.

Pentru înmulțirea clasică a două matrice se folosește operatorul '%*%' aplicat unor matrice de dimensiuni corespunzătoare.

```
> C=matrix(c(2,4,1,10),ncol=2,byrow=TRUE); C
      [,1] [,2]
[1,]     2     4
[2,]     1    10

> D=matrix(c(4,6,8,5),ncol=2,byrow=TRUE); D
      [,1] [,2]
[1,]     4     6
[2,]     8     5
> C%*%D      #inmultirea clasica dintre doua matrice
      [,1] [,2]
[1,]    40    32
[2,]    84    56
```

Dacă unul dintre operanzi este vector și al doilea matrice, atunci la înmulțirea vectorului cu matricea, vectorul este organizat în vector linie sau coloană astfel încât să se poată realiza înmulțirea dintre vector și matrice. De exemplu:

```
> x=1:2; x      #vector
[1] 1 2
> A=matrix(1:4,c(2,2)); A      #matrice
      [,1] [,2]
[1,]     1     3
[2,]     2     4
> x%*%A
      [,1] [,2]
[1,]     5    11
> A%*%x
      [,1]
[1,]     7
[2,]    10
```

$$(1 \ 2) \cdot \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} = (5 \ 11)$$

$$\begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 7 \\ 10 \end{pmatrix}$$

Pentru determina transpusa unei matrice se poate folosi funcția t:

```
> A=matrix(c(2,4,1,10,6,0),ncol=2,byrow=TRUE); A
> A
      [,1] [,2]
[1,]     2     4
[2,]     1    10
[3,]     6     0
> t(A)      # transpusa matricei A
      [,1] [,2] [,3]
[1,]     2     1     6
[2,]     4    10     0
```

Pentru a determina inversa unei matrice se poate folosi funcția solve:

```
> M=matrix(c(1,2,3,4),ncol=2,byrow=TRUE); M
      [,1] [,2]
[1,]     1     2
[2,]     3     4
```

```

> solve(M)      # M^(-1)
      [,1] [,2]
[1,] -2.0  1.0
[2,]  1.5 -0.5

> A=matrix(c(1,2,1,2),c(2,2)); A      #det(A)=0
      [,1] [,2]
[1,]     1     1
[2,]     2     2
> solve(A)
Error in solve.default(A) :
  Lapack routine dgesv: system is exactly singular

```

2.4. Tablouri numerice

Tablourile de numere (*arrays*) sunt generalizări ale matricelor la dimensiuni mai mari sau egale cu doi. O matrice este un *array* de dimensiune doi.

Crearea unui tablou de numere se poate face utilizând funcția **array(elemente,dimensiuni)**. De exemplu:

```

> m5=array(1:12, dim=c(2,3,2))
> m5
, , 1
      [,1] [,2] [,3]
[1,]     1     3     5
[2,]     2     4     6

, , 2
      [,1] [,2] [,3]
[1,]     7     9    11
[2,]     8    10    12

```

O altă modalitate de a crea tablouri multidimensionale, este de a utiliza funcția **dim** ca în exemplul de mai jos:

```

> x=1:8
> dim(x)=c(2,2,2)
> x
, , 1
      [,1] [,2]
[1,]     1     3
[2,]     2     4

, , 2
      [,1] [,2]
[1,]     5     7
[2,]     6     8

```

Astfel, vectorul *x* a fost transformat în tablou de dimensiuni $2 \times 2 \times 2$.

Tablourile pot fi utilizate în diferite operații aritmetice, multe dintre acestea fiind similare celor descrise pentru matrice. Se pot efectua de asemenea operații între tablouri numerice și alte obiecte R precum ar fi vectorii și matricele. De exemplu:

```

> x=1:4
> y=5:8
> t1=array(x,c(2,2,2))
> t2=array(y,c(2,2,2))
> t1
, , 1
    [,1] [,2]
[1,]    1    3
[2,]    2    4

, , 2
    [,1] [,2]
[1,]    1    3
[2,]    2    4

> t2
, , 1
    [,1] [,2]
[1,]    5    7
[2,]    6    8

, , 2
    [,1] [,2]
[1,]    5    7
[2,]    6    8

> t1*t2
, , 1
    [,1] [,2]
[1,]    5   21
[2,]   12   32

, , 2
    [,1] [,2]
[1,]    5   21
[2,]   12   32

> t1%%t2 #suma_k suma_i,j t1[i,j]*t2[i,j]
    [,1]
[1,]  140

```

Dacă A și B sunt două tablouri numerice atunci **produsul lor exterior** (engl., *outer product*) va fi un tablou numeric ale cărui elemente se obțin prin formarea tuturor combinațiilor posibile ale elementelor lui A și B . Simbolul pentru acest tip de produs este `'%o%'`:

```
> AB = A%o%B
```

sau echivalent

```
> AB = outer(A,B,"*")
```

De remarcat că produsul definit de simbolul `'%*%'` (înmulțirea clasică dintre matrice) este diferit de produsul reprezentat prin simbolul `'%o%'`.

Funcția `aperm()` poate fi utilizată pentru a permuta dimensiunile unui tablou numeric. Al doilea argument este o permutare de numere întregi $\{1, 2, \dots, k\}$, unde k este numărul de indici al tabloului numeric. Rezultatul acestei funcții este un tablou numeric de aceeași dimensiuni cu dimensiunile, respectiv datele permutate conform

celui de-al doilea argument. Astfel, funcția `aperm` poate fi o transpusă generalizată ca în exemplul următor:

```
A=array(1:12, dim=c(2,3,2))
> A
, , 1
      [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6

, , 2
      [,1] [,2] [,3]
[1,]    7    9   11
[2,]    8   10   12

> B = aperm(A,c(2,1,3))
> B
, , 1
      [,1] [,2]
[1,]    1    2
[2,]    3    4
[3,]    5    6

, , 2
      [,1] [,2]
[1,]    7    8
[2,]    9   10
[3,]   11   12
```

În cazul particular al unei matrice `A`, instrucțiunea `B=aperm(A,c(2,1))` este echivalentă cu instrucțiunea `B=t(A)` și determină transpusa matricei `A`.

2.5. Factori

Conceptual, factorii în R sunt variabile care iau un număr limitat de valori diferite. Aceștia sunt reprezentanți intern ca valori întregi $1, 2, \dots, k$, unde k este numărul de niveluri. Variabilele de tip numeric sau cele de tip caracter pot fi convertite în factori. Nivelurile de factori sunt întotdeauna șiruri de caractere. Factorii pot fi utilizați pentru sortarea elementelor vectorului, cât și pentru ordonarea după o anumită funcție de distribuție.

Pentru crearea unui obiect de tip factor se folosește funcția `factor`. Să considerăm un exemplu în care avem 10 studenți cazați la un cămin studentesc, aceștia fiind identificați prin identificatorul județului din care provin:

```
> orase=c("VL", "DB", "VL", "AG", "OT", "SB", "DB", "VL", "AG", "AG")
> orase
[1] "VL" "DB" "VL" "AG" "OT" "SB" "DB" "VL" "AG" "AG"
> orasef=factor(orase)           #generarea obiectului factor
> orasef
[1] VL DB VL AG OT SB DB VL AG AG
Levels: AG DB OT SB VL
> levels(orasef)                #lista judetelor in ordine alfabetica
[1] "AG" "DB" "OT" "SB" "VL"
```


Ordonarea șirurilor de caractere în ordine alfabetică se poate face și cu funcția `sort`:

```
> sort(orsef)
[1] AG AG AG DB DB OT SB VL VL VL
Levels: AG DB OT SB VL
```

Continuând aplicația anterioară, să presupunem că mai știm mediile generale de admitere, respectiv mediile de pe anul universitar anterior:

```
> medii = c(7.9, 9, 8.5, 7.7, 8.6, 9.3, 7.3, 9.5, 8.8, 9.2)
```

Să presupunem că dorim să știm media generală medie și erorile standard corespunzătoare fiecărui oraș în parte. Pentru ceasta putem folosi funcția `tapply`:

```
> media_medie=tapply(medii,orsef,mean)
> media_medie
      AG      DB      OT      SB      VL
8.566667 8.150000 8.600000 9.300000 8.633333
```

Aici, `tapply` folosește funcția `mean` pentru a grupa elementele primului argument, *medii*, definită de nivelurile celei de-a doua componente, *orsef*, ca și cum ar fi structuri vectoriale separate. Rezultatul este un vector de lungime egală cu numărul nivelurilor. Pentru a calcula și eroarea standard, determinăm variația cu funcția `var` și calculăm eroarea standard:

```
> eroare_std=function(x){sqrt(var(x))/length(x)}
> media_medie=tapply(medii,orsef,eroare_std)
> media_medie
      AG      DB      OT      SB      VL
0.2589151 0.6010408 NA NA 0.2694301 #eroarea std
```

Valorile NA corespund cazului în care calculul erorii standard nu a fost posibil, deoarece numărul studenților a fost mai mic sau egal cu 1.

Factorii ordonați se obțin practic prin ordonarea nivelurilor. Pentru a crea un factor ordonat sau pentru a ordona un factor existent se poate utiliza funcția `ordered`. Nivelurile unui factor ordonat specifică poziția acestora pe o scară ordinară. Pe lângă factorii ordonați, mai există și factorii de contrast.

2.6. Data frame

`Data frame` este o colecție de date organizate într-un tablou dreptunghiular în care pe fiecare coloană sunt date de același tip. Un exemplu este tabelul următor:

Nume	Analiză	Algebră
Popescu Nicolae	8	6
Stan Ion	7	7
Vasile Anca	9	8

O variabilă `data frame` se creează prin alipirea unor vectori de aceeași dimensiune, ca în exemplele de mai jos:

```
> c1=c("a","b","c")
> c2=c(1,2,3)
> df1=data.frame(c1,c2)      #create data frame
> df1
```

```

      c1 c2
1     a  1
2     b  2
3     c  3

```

sau

```

> clasa=c("9A","9B","9C","9D","9E","9F","9G","9H")
> media=(8.82,9.01,8.45,8/23,8.14,8.01,7.76,9.08)
> catalog=data.frame(clasa,media); catalog
  clasa      media
1    9A  8.820000
2    9B  9.010000
3    9C  8.450000
4    9D  0.3478261
5    9E  8.140000
6    9F  8.010000
7    9G  7.760000
8    9H  9.080000

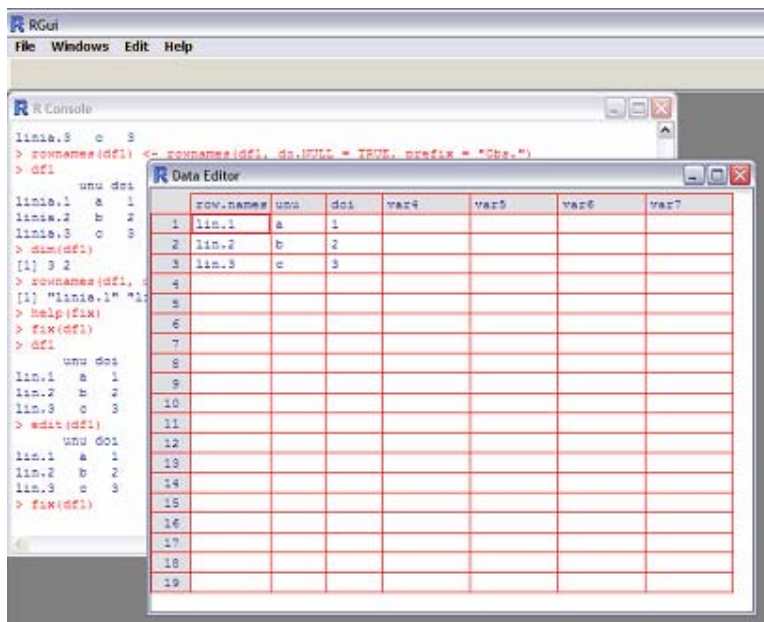
```

Numele coloanelor coincide cu numele vectorilor, iar liniile au numele 1, 2, 3, ...
Putem schimba aceste nume ca în secvența de mai jos:

```

> rownames(df1)=c("linia.1","linia.2","linia.3")
> colnames(df1)=c("unu","doi")
> df1
      unu doi
linia.1  a  1
linia.2  b  2
linia.3  c  3

```



Numele liniilor, coloanelor ca și valorile din data.frame (sau din alt obiect precum matricea sau vectorul) pot fi editate cu ajutorul unui editor propriu R care se apelează prin comenzile **fix** sau **edit**, ca mai jos:

```

> fix(df1)
> df1
      unu doi
lin.1  a  1
lin.2  b  2
lin.3  c  3

```

În procesul de editare am schimbat numele liniilor.

Se pot crea matrice sau data frames cu ajutorul editorului încorporat prin instrucțiunile

```

> df2=edit(data.frame())
> m5=edit(matrix())

```

Valorile pentru df2 și m5 vor fi cele completate în editor.

Numele coloanelor unei date de tip `data frame` se obține prin comanda `names(data.frame)`, iar coloana cu un nume dat se poate accesa prin `data$nume` ca în exemplele următoare:

```
> df1
      unu doi
lin.1   a   1
lin.2   b   2
lin.3   c   3
> names(df1)      #numele coloanelor
[1] "unu" "doi"
> df1$unu          #coloana cu numele "unu"
[1] a b c
Levels: a b c
> df1$unu[3]       #elem. al 3-lea de pe coloana "unu"
[1] c
Levels: a b c
```

sau

```
> names(catalog)      #numele coloanelor
[1] "clasa" "media"
> catalog[1]          #prima coloana
      clasa
1      9A
2      9B
3      9C
4      9D
5      9E
6      9F
7      9G
8      9H
> catalog[1,1]        #elem.din prim coloana, prima linie
[1] 9A
Levels: 9A 9B 9C 9D 9E 9F 9G 9H
```

O variabilă de tip `data.frame` este de fapt un tabel. Poate fi gândită ca o matrice care are un cap de tabel în care și liniile și coloanele au denumiri. Liniile sunt considerate „nivele”, iar coloanele sunt vectori de același tip care poarta un nume.

Datele `data frame` pot fi introduse în R și din fișiere text exterioare prin comanda **`read.table`**. Datele în fișier trebuie să umple o zonă dreptunghiulară. Fișierul text cu date poate fi copiat în clipboard, poate fi pe un suport de memorie, într-o locație de internet. Din mulțimea de parametri ai instrucțiunii `read.table` este important de reținut: numele fișierului, `sep`=semnul separator între date, `header` care poate fi `TRUE` dacă primul rând este citit ca nume ale variabilelor, altfel `FALSE`, `dec` ce indică semnul separator pentru zecimale (implicit). Liniile vide ca și cele care încep cu semnul `#` sunt ignorate. Datele citite sunt interpretate ca `data frame`. O comandă asemănătoare este **`read.csv`**. În acest caz fișierul trebuie să fie un fișier text cu datele separate prin virgulă.

	1	2	3	4	5	6	7
1							
2		Analiza	Algebra	Fizica			
3	Ion	8	7	7			
4	Nicu	6	5	8			
5	Vasile	8	9	7			
6							
7							
8							
9							

De exemplu, să presupunem că s-au creat niște date în Excel ca în figura alăturată, se selectează și se copiază în clipboard. Datele se citesc în R astfel:

```
> t=read.table("clipboard")
> t
```

	Analiza	Algebra	Fizica
Ion	8	7	7
Nicu	6	5	8
Vasile	8	9	7

Pentru a vedea ce tip de dată este *t* apelăm funcția `class`:

```
> class(t)
[1] "data.frame"
```

Dacă am fi selectat doar matricea de date numerice s-ar fi citit în R tot ca *data frame*. O metodă alternativă este să salvăm documentul Excel ca fișier text sau fișier „csv” și să-l citim cu `read.table` sau `read.csv`.

```
> t2=read.table(file.choose())
> t2
```

	Analiza	Algebra	Fizica
Ion	8	7	7
Nicu	6	5	8
Vasile	8	9	7

```
> t3=read.csv(file.choose())
> t3
```

	X	Analiza	Algebra	Fizica
1	Ion	8	7	7
2	Nicu	6	5	8
3	Vasile	8	9	7

Funcția **file.choose** deschide o fereastră de dialog pentru alegerea fișierului. Dacă se pune numele fișierului atunci fișierul trebuie să fie în directorul curent sau să fie indicat cu cale cu tot. În rezumat avem:

```
t=read.table("clipboard")
#pentru citirea din clipboard
t=read.table(file.choose())
#citire dintr-un fișier text
t=read.table("numele.fișierului")
#citire dintr-un fișier aflat în directorul curent
t=read.table("http://adresa.web/fisier.txt")
#citire fișier de pe internet.
```

Pentru a scrie un obiect *data frame* într-un fișier se folosește instrucțiunea **write.table** sau **write.csv** care au cam aceeași parametri ca `read.table`. De exemplu prin instrucțiunea

```
> write.table(t, file="t.txt")
```

se salvează obiectul data frame `t` în fișierul "t.txt" din directorul curent de unde poate fi citit eventual altă dată cu `read.table`.

Pentru scrierea vectorilor sau matricelor în fișiere text se poate utiliza instrucțiunea `write`. Citirea matricelor din aceste fișiere se face cu instrucțiunea `scan`. Sintaxa instrucțiunii `write` este:

```
write(matricea, file = "fișierul", ncolumns  
=nr.coloane, append = FALSE, sep = ",")
```

din care doar matricea (vectorul) și fișierul sunt obligatorii. Dacă `file=""` atunci se scrie la consolă. Dacă `append=TRUE` atunci se adaugă datele la sfârșitul fișierului altfel se șterge fișierul și se recrează cu datele noi. Separatorul poate fi și alt caracter, de exemplu spațiu sau tab. Scrierea se face pe linii. Pentru a scrie coloanele la început, trebuie transpusă matricea. Citirea unei matrice dintr-un fișier text cu instrucțiunea `scan` are forma

```
scan(file = "", what = numeric), nmax = -1, n = -1,  
sep = "", dec = ".", skip = 0, nlines = 0)
```

`File` este numele fișierului, `what` reprezintă tipul de date care se citesc (logical, integer, numeric, complex, character, raw, list), `nmax` reprezintă numărul maxim de date ce se citesc (dacă lipsește sau nu e pozitiv atunci se citește până la sfârșitul fișierului), `n` are aceeași semnificație doar că se ignoră valorile invalide, `sep` este separatorul de date (dacă lipsește se presupune că datele sunt separate prin spații), `dec` este separatorul pentru zecimale în numerele reale (implicit), `skip` reprezintă numărul de linii de la început care sunt ignorate, `nlines` este numărul maxim de linii care se citesc.

R include și un număr de **baze de date** care se pot folosi în diverse aplicații. Aceste baze de date sunt de fapt data frame-uri.

Pentru a vizualiza toate bazele de date disponibile, se poate folosi funcția `data`:

```
> data()
```

și se va deschide o fereastră în care sunt enumerate toate bazele de date disponibile. Pentru a afla informații suplimentare despre o anumită bază de date, se poate apela help-ul sistemului R:

```
> ?cars
```

Se va deschide o fereastră web în care se specifică faptul că baza de date conține 50 de observații luate din anul 1920, iar viteza este numerică și măsurată în mph, iar distanța este de asemenea numerică, măsurată în ft.

Pentru a accesa una dintre aceste baze de date, se poate scrie `data(numele bazei de date)`. De exemplu, pentru a accesa baza de date `cars` care conține vitezele și distanțele de oprire ale unor mașini putem proceda astfel:

```
> data(cars)  
> cars[1:5,] #primele 5 rânduri ale bazei de date  
  speed dist #coloanele bazei de date  
1      4    2  
2      4   10  
3      7    4  
4      7   22  
5      8   16
```

Dacă se dorește manipularea pe caracteristici, atunci se poate folosi vectorul **`cars[,1]`** ce ne dă toate vitezele măsurate sau vectorul **`cars[,2]`** ce ne dă toate distanțele disponibile.

O altă modalitate de a lucra cu o bază de date este de a folosi funcția **`attach`**:

```
> attach(cars)
> names(cars)      #numele coloanelor bazei de date
[1] "speed" "dist"
> mean(speed)      #media datelor din prima coloanei
[1] 15.4
> mean(cars[,1])   #o varianta echivalenta
[1] 15.4
> mean(cars$speed) #a doua varianta echivalenta
[1] 15.4
```

2.7. Liste

Lista este o colecție de obiecte R care pot fi de tipuri diferite și de mărimi diferite. Crearea unei liste se face simplu prin funcția **`list`** urmată între paranteze de numele obiectelor ce o compun. Să presupunem că există o familie cu numele tatălui „Ion”, numele mamei „Maria”, iar numele copiilor sunt: „Andrei” (4 ani), „Alexandra” (7 ani) și „Bogdan” (9 ani).

```
> lista1=list(sot="Ion",sotie="Maria",nr.copii=3,
              varsta.copii=c(4,7,9))
> lista1      #afisarea continutului listei
$sot
[1] "Ion"
$sotie
[1] "Maria"
$nr.copii
[1] 3
$varsta.copii
[1] 4 7 9
```

iar accesul la un element din listă se face prin **`numele_listei[[numarul de ordine]]`** al elementului ca în comanda:

```
> lista1[[2]]
[1] "Maria"
```

Alt exemplu:

```
> clasa=c("9A","9B","9C","9D","9E","9F","9G","9H")
> media=c(8.82,9.01,8.45,8/23,8.14,8.01,7.76,9.08)
> catalog=data.frame(clasa,media)
> lista=list(clasa,media,catalog)
> lista[1]
[[1]]
[1] "9A" "9B" "9C" "9D" "9E" "9F" "9G" "9H"
> lista[2]
[[1]]
[1] 8.8200000 9.0100000 8.4500000 9.3478261
[5] 8.1400000 8.0100000 7.7600000 9.0800000
> lista[3]
[[1]]
```

```

      clasa      media
1      9A 8.8200000
2      9B 9.0100000
3      9C 8.4500000
4      9D 9.3478261
5      9E 8.1400000
6      9F 8.0100000
7      9G 7.7600000
8      9H 9.0800000
> lista[[1]][2]
[1] "9B"

```

2.8. Raw-uri

Șirurile de biți se obțin prin declararea unei variabile ca vector de tip `raw` și atribuirea fiecărui element din vector a unei valori întregi între 0 și 255 convertită la tipul `raw` prin instrucțiunea `as.raw`. Dacă nu reușește conversia atunci valoarea este luată implicit 0. Exemplu:

```

> x=raw(2) #creeaza un vector raw gol de dim. 2
> x[1]=as.raw(100)      #codul ASCII 100 (caract. d)
> x[2]=charToRaw("Z") #codul ASCII al caract. Z
> x #afisare a cate o pereche de cifre hexazecimale
[1] 64 5a
> rawToChar(x)  #conversie raw in sir de caractere
[1] "dZ"
> rawToBits(x) #conversie la biti
[1] 00 00 01 00 00 01 01 00 00 01 00 01 01 00 01 00

```

Vectorii de tip `raw` pot fi manipulați folosind operatorii la nivel de biți: `!` (complement față de 1), `&` (și), `|` (sau) and `xor` (sau exclusiv). Continuând exemplul de mai sus:

```

> y=!x          #complement data de 1
> rawToBits(y)
[1] 01 01 00 01 01 00 00 01 01 00 01 00 00 01 00 01
> z=charToRaw("AB"); z
[1] 41 42
> rawToBits(z)
[1] 01 00 00 00 00 00 01 00 00 01 00 00 00 00 01 00
> rawToBits(xr&z)      #si la nivel de biti
[1] 00 00 00 00 00 00 01 00 00 01 00 00 00 00 01 00
> rawToBits(xr|z)      #sau la nivel de biti
[1] 01 00 01 00 00 01 01 00 00 01 00 01 01 00 01 00

```

Capitolul 3

Funcții predefinite în R

Practic totul în R se obține prin funcții. Să urmărim comenzile de mai jos:

```
> "+"(2, 3)
[1] 5
> 2+3
[1] 5
```

Se vede că operația de adunare este implementată ca o funcție. La fel alte operații $x-y$, x/y , $x*y$, x^y sunt funcții pentru care utilizăm forma operatorială de notație, mai comodă.

R vine cu un mare număr de funcții predefinite pentru prelucrarea datelor, iar diversele pachete adiționale vin cu funcții suplimentare. În cele ce urmează descriem un număr de funcții curent utilizate. În manualele ce însoțesc distribuțiile de R, în cărțile din bibliografie sau în diverse locații de pe internet (de exemplu la adresa <http://www.statmethods.net/> sau pur și simplu căutând cu Google documentație pentru R) se pot găsi informații pentru multe alte funcții disponibile. Pe lângă funcțiile predefinite putem defini funcții noi care se pot utiliza exact ca cele predefinite.

Un lucru este important de menționat. Argumentele funcțiilor au în general nume și dacă se specifică numele lor atunci nu este importantă ordinea în care se pun, ca în exemplul următor în care se crează un șir de numere cu funcția `seq`.

```
> seq(from=1, by=2, to=5)
[1] 1 3 5
> seq(by=2, to=5, from=1)
[1] 1 3 5
```

De asemenea, unele funcții admit un număr mai mare de argumente care nu sunt compatibile decât în anumite combinații, ca în exemplul de mai jos :

```
> seq(by=2, to=5, from=1)
[1] 1 3 5
> seq(from=1, by=2, length=3)
[1] 1 3 5
> seq(from=1, by=2, to=5, length=3)
Error in seq.default(from = 1, by = 2, to = 5,
length = 3): too many arguments
```

Numele argumentelor formale este ales intuitiv (în engleză) și după un pic de practică nu este dificil a le ține minte. De la o funcție la alta unele argumente formale au valori predefinite și funcțiile pot fi apelate fără aceste argumente (se utilizează valorile implicite în calcul). Apar astfel situații când aceeași funcție este apelată cu un număr diferit de argumente. O astfel de situație este curentă pentru funcțiile grafice de exemplu. Dacă funcțiile sunt apelate fără menționarea numelui argumentelor atunci ordinea și semnificația argumentelor este strictă, ca în programul de definiție a funcției. De exemplu

```
> seq(1, 5, 2)
[1] 1 3 5
```


pune în evidență că în funcția `seq` primul argument este `from`, al doilea este `to` și al treilea este `by`. Pentru foarte multe funcții din R apelul se face cu primele 1, 2, 3 argumente nedenumite, dar în ordinea și cu semnificația specificată în documentație, pe când restul argumentelor sunt parametri opționali din care unii se introduc prin numele lor iar ceilalți se utilizează prin valorile lor implicite.

3.1. Funcții numerice

Aceste funcții aplicate unui număr dau ca rezultat un alt număr, iar aplicate unui vector sau unei matrice se aplică fiecărui element al vectorului sau matricei respective (sunt vectorizate).

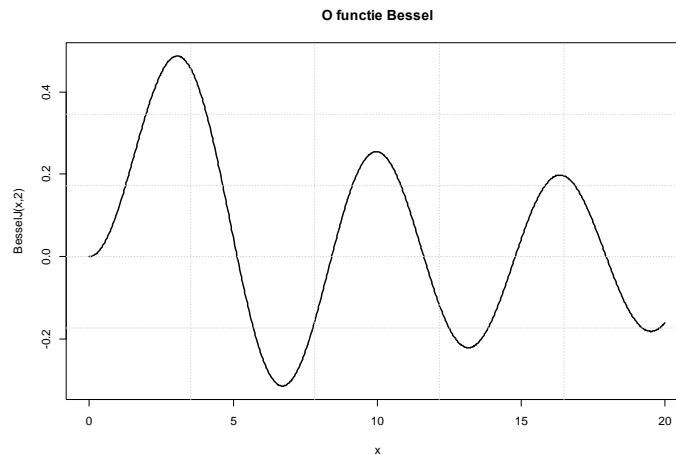
Funcția	Explicații
<code>abs(x)</code>	$ x $
<code>sqrt(x)</code>	\sqrt{x}
<code>ceiling(x)</code>	Dacă $x \in (n-1, n]$ valoarea returnată este n
<code>floor(x)</code> , <code>trunc(x)</code>	Partea întreagă a lui x
<code>round(x, digits=n)</code>	Rotunjire la n zecimale
<code>signif(x, digits=n)</code>	Rotunjire la cele mai semnificative n cifre
<code>cos(x)</code> , <code>sin(x)</code> , <code>tan(x)</code> , <code>acos(x)</code> , <code>asin(x)</code> , <code>atan(x)</code> , <code>cosh(x)</code> , <code>acosh(x)</code> , etc.	Funcții trigonometrice directe și inverse precum și funcții hiperbolice directe și inverse
<code>log(x)</code>	Logaritmul natural, $\ln x$
<code>log10(x)</code>	Logaritmul în baza 10, $\lg x$
<code>exp(x)</code>	e^x
<code>besselI(x, nu)</code> <code>besselK(x, nu)</code> <code>besselJ(x, nu)</code> <code>besselY(x, nu)</code>	Funcții Bessel. Parametru nu este ordinul funcției. Parametru x este un vector numeric cu valori pozitive.
<code>gamma(x)</code> , <code>lgamma(x)</code>	Funcția gamma și logaritmul natural din gamma
<code>beta(a,b)</code> , <code>lbeta(a,b)</code>	Funcția beta și logaritmul natural din beta
<code>digamma(x)</code>	Calculează derivata lui $\log(\text{gamma}(x))$
<code>psigamma(x, deriv)</code>	Calculează derivata de ordin $deriv$ a lui $\text{digamma}(x)$
<code>choose(n,k)</code> , <code>lchoose(n,k)</code>	Combinări de n câte k (n poate fi real); <code>lchoose</code> calculează logaritmul natural din <code>choose</code> .
<code>factorial(x)</code>	Calculează $\text{gamma}(x+1)$

<code>F<-splinefun(x, y, method)</code>	Valoarea întoarsă este o funcție, F , care poate fi ulterior apelată normal, $F(a)$ pentru valoarea funcției sau $F(2, deriv=1)$ ceea ce duce la calculul derivatei funcției spline. <code>deriv</code> poate lua valorile 0, 1, 2, 3. x, y sunt vectorii cu coordonatele punctelor de interpolat <code>method = "fmm", "periodic", "natural", "monoH.FC"</code> , și indică tipul de funcție spline
<code>spline(x, y, n, method, xmin, xmax, xout)</code>	Forma aceasta calculează valoarea funcției spline determinată de vectorii x și y în n puncte echidistante între $xmin$ și $xmax$ sau în vectorul $xout$.

Exemplul 3.1 Să se reprezinte grafic funcția Bessel J_2 pentru $x \in [0,20]$.
O posibilă rezolvare este:

```
x=seq(0,20,length=1000)
y=besselJ(x,2)
plot(x,y,type="lines",lwd=2,ylab="BesselJ(x,2)",
xlab="x", main="O functie Bessel")
grid(5,5)
```

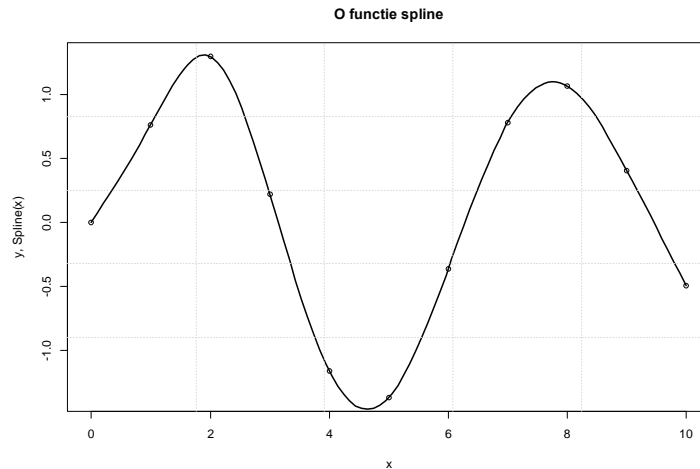
Graficul este :



Exemplul 3.2 Se dau doi vectori și să se reprezinte grafic punctele care au coordonatele date de cei doi vectori precum și funcția spline naturală ce trece prin puncte. O posibilă rezolvare este:

```
x=0:10
y=x*sin(x)/(1+x^2/10)
plot(x,y,main="O functie spline", xlab="x", ylab="y",
Spline(x) ")
F<-splinefun(x,y,method="natural")
x1=seq(0,10,length=100)
y1=F(x1)
lines(x1,y1,lwd=2)
grid(5,5)
```

Graficul este:



3.2. Funcții pentru matrice

Matricele reprezintă un obiect matematic important în statistică și de aceea R are implementate multe funcții referitoare la ele. Mai jos indicăm unele din ele.

Operatorul sau funcția	Explicații
<code>A * B</code>	Înmulțirea element cu element
<code>A %*% B</code>	Înmulțirea matriceală
<code>outer(a,b,f)</code> <code>a %o% b</code>	Dacă a și b sunt vectori generează o matrice x cu $x_{i,j} = f(a_i, b_j)$. Dacă a și b sunt matrice generează un array x cu $x_{i,j,k,l} = f(a_{i,j}, b_{k,l})$. Forma <code>a %o% b</code> este asemănătoare față de funcția <code>produs</code> .
<code>crossprod(A,B)</code> <code>crossprod(A)</code>	$A'B$, respectiv $A'A$
<code>cbind(A,B,...)</code>	Combină orizontal (rezultatul va avea liniile din A, B, \dots)
<code>rbind(A,B,...)</code>	Combină vertical (rezultatul este o matrice cu coloanele din A, B, \dots)
<code>rowMeans(A)</code>	Se obține un vector cu mediile liniilor
<code>rowSums(A)</code>	Se obține un vector cu sumele elementelor liniilor
<code>colMeans(A)</code>	Se obține un vector cu mediile coloanelor
<code>colSums(A)</code>	Se obține un vector cu sumele elementelor coloanelor
<code>min(A)</code> , <code>max(A)</code> , <code>range(A)</code>	Minimul, maximul, respectiv un vector cu minimul și maximul
<code>solve(A)</code>	Inversa matricei A
<code>ginv(A)</code>	Inversa generalizată Moore-Penrouse (se găsește în pachetul MASS)
<code>x<-eigen(A)</code>	<code>x\$values</code> conține valorile proprii ale lui A <code>x\$vectors</code> conține vectorii proprii ai lui A

<code>t(A), det(A)</code>	Transpusa, respectiv determinantul unei matrice
<code>x<-svd(A)</code>	Descompunerea singulară a matricei A , $A=UDV^t$, cu $D=\text{diag}(d)$ $x\$d$ = vector ce conține valorile singulare ale matricei A $x\$u$ = matrice ortogonală cu vectorii singulari la stânga ai matricei A $x\$v$ = matrice ortogonală cu vectorii singulari la dreapta ai matricei A
<code>R <- chol(A)</code>	Factorizarea Choleski, $A=R^tR$, cu R superior triunghiulară, iar A pozitiv definită.
<code>x<-qr(A)</code>	Descompunerea QR a matricei A $x\$qr$ are deasupra diagonalei principale matricea R (superior triunghiulară), iar sub diagonala principală informații pentru construcția matricei Q . $x\$rank$ este rangul lui A $x\$qraux$ este un vector cu informații suplimentare despre Q $x\$pivot$ conține informații despre strategia utilizată
<code>qr.Q(x), qr.R(x), qr.X(x)</code>	Crează dintr-un obiect x de tip qr generat de funcția <code>qr</code> matricile Q , R sau matricea inițială x .
<code>v<-sort(A,...)</code>	$v\$x$ conține elementele vectorului sau matricei A (matricea este convertită în vector punând coloanele una după alta) sortate. $v\$ix$ conține dacă parametrul suplimentar <code>index.return=FALSE</code> un vector cu pozițiile originale ale elementelor ordonate. Dacă <code>decreasing=TRUE</code> sortarea se face descrescător Dacă <code>partial=vector</code> de întregi sortarea se face numai pentru componentele indicate de <code>partial</code>
<code>rev(x)</code>	Inversează ordinea elementelor vectorului x

Exemplul 3.3 Se generează aleator un sistem liniar și se determină soluția lui.

```
a = matrix(runif(25),ncol=5)
b = runif(5)
x = solve(a,b)
a
b
x
```

Se obține răspunsul:

```
> a
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 0.2368829 0.79721628 0.1265954 0.04369927 0.2061009
[2,] 0.1015411 0.70852984 0.7379991 0.03374915 0.4783938
[3,] 0.4941975 0.28829965 0.2681115 0.92383512 0.3105410
[4,] 0.8037851 0.04331768 0.5082257 0.62702404 0.3799365
[5,] 0.2050783 0.22855288 0.3172556 0.75981321 0.1595518
```

```
> b
[1] 0.4920966 0.5553061 0.3207777 0.2010808 0.5628473
> x
[1] 0.1839354 0.9262001 1.2986354 0.3486333 -2.2779749
```

Exemplul 3.4 Să se determine în cât timp se rezolvă un sistem liniar 1000×1000 generat aleator. Rezolvarea este în secvența de cod următoare:

```
ptm <- proc.time()
a = matrix(runif(1000000), ncol=1000)
b = runif(1000)
x = solve(a,b)
proc.time() - ptm
```

Răspunsul pe un calculator cu procesor Q9300 quad la 2.5 Ghz este:

```
user system elapsed
0.69 0.08 1.36
```

Exemplul 3.5 Să se determine valorile și vectorii proprii ai unei matrice precum și descompunerea ei singulară. Se va alege o matrice singulară.

```
x=runif(5); y=runif(5)
a=outer(x,y,"+")
v<-eigen(a)
ds<-svd(a)
```

Ca produs `outer`, matricea a este singulară. Valorile și vectorii proprii sunt:

```
> v
$values
[1] 6.819055e+00 -7.174162e-02 -2.935982e-16 -5.941790e-17
-1.790505e-17

$vectors
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 0.3264044 0.73432693 -0.67659694 0.3795791 0.3029054
[2,] 0.4377868 0.04995971 0.68624453 -0.3881873 -0.1576330
[3,] 0.5506698 -0.64362770 0.21134194 0.6398041 0.3044435
[4,] 0.4216151 0.14932348 -0.14367175 -0.5354580 -0.8120513
[5,] 0.4699004 -0.14735628 -0.07731778 -0.0957380 0.3623354
```

Se vede că trei valori proprii sunt aproape zero. În fapt ele sunt zero, matricea având rangul 1 și doar erorile de rotunjire fac rezultatele de ordinul 10^{-16} . Descompunerea singulară este:

```
> ds
$d
[1] 6.870905e+00 1.077193e-01 1.869400e-16 2.686574e-17 3.540473e-18

$u
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] -0.3262407 0.77055695 -0.5096937 -0.1995559 0.014098365
[2,] -0.4377668 0.09392287 0.3454248 0.1386439 -0.813017974
[3,] -0.5507955 -0.59182723 -0.3600220 -0.4655446 -0.004147016
[4,] -0.4215743 0.19216386 0.6827627 -0.2779026 0.491887645
[5,] -0.4699219 -0.10116352 -0.1584718 0.8043592 0.311178980
```

```

$V
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] -0.3932570 -0.5884878  0.674404332 -0.20995153 -0.01141111
[2,] -0.3959388 -0.5602933 -0.706356531  0.03387009  0.17093016
[3,] -0.4729784  0.2496549 -0.138133931 -0.21376690 -0.80572118
[4,] -0.4682765  0.2002217  0.164709931  0.84032560  0.08574326
[5,] -0.4955703  0.4871727  0.005376199 -0.45047726  0.56045887

```

Putem verifica faptul că $A=UDV^t$, unde U și V sunt matrici ortogonale (date de `ds$u`, `ds$V`), iar D este matricea diagonală cu `ds$d` pe diagonală. Verificarea este făcută prin comanda următoare:

```

> a-ds$u%*%diag(ds$d)%*%t(ds$V)
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] -1.110223e-16 -8.881784e-16 -8.881784e-16 -6.661338e-16 -8.881784e-16
[2,]  0.000000e+00 -4.440892e-16 -2.220446e-16 -2.220446e-16 -2.220446e-16
[3,] -2.220446e-16 -4.440892e-16 -2.220446e-16 -2.220446e-16 -4.440892e-16
[4,] -2.220446e-16 -4.440892e-16 -2.220446e-16 -2.220446e-16 -2.220446e-16
[5,] -2.220446e-16 -6.661338e-16 -4.440892e-16 -4.440892e-16 -4.440892e-16

```

Din nou erorile de rotunjire duc la un rezultat doar foarte apropiat de zero, nu exact zero.

3.3. Funcții pentru șiruri de caractere

Pentru prelucrarea șirurilor și vectorilor cu elemente șiruri de caractere se pun la dispoziție mai multe funcții. Unele sunt în tabelul de mai jos.

Funcția	Explicații
<code>character(length)</code>	Crează un șir de blank-uri de lungime <code>length</code>
<code>nchar(x,...)</code>	Numărul de caractere din șir
<code>substr(x,start=n1,stop=n2)</code>	Extrage un subșir dintr-un șir dat sau înlocuiește o parte dintr-un subșir cu un șir dat
<code>grep(model,x,fixed=TRUE)</code>	Caută șirul <code>model</code> în <code>x</code> (<code>x</code> este convertit la un vector de șiruri). Rezultatul este o mulțime de indici care specifică indicii șirurilor unde s-a găsit <code>model</code> . Printre parametrii opționali se numără: <code>fixed=TRUE, FALSE</code> care dacă este <code>FALSE</code> căutarea se face după <code>model=expresie regulată</code> , <code>ignore.case=TRUE, FALSE</code> , care specifică dacă se ignoră sau nu faptul că sunt litere mari sau mici.
<code>sub(sir1,sir2,x,fixed=TRUE)</code> <code>gsub(sir1,sir2,x,fixed=TRUE)</code>	Înlocuiește <code>sir1</code> cu <code>sir2</code> în <code>x</code> . Dacă <code>fixed=FALSE</code> atunci <code>sir1</code> este o expresie regulată. Parametrul <code>ignore.case=TRUE, FALSE</code> controlează dacă se ține seama de caractere majuscule și minuscule. Funcția <code>sub</code> înlocuiește doar prima apariție a lui <code>sir1</code> iar funcția <code>gsub</code> înlocuiește toate aparițiile șirului <code>sir1</code> .

<code>strsplit(x, split)</code>	Scindează şirul <i>x</i> în subşiruri. Scindarea se face acolo unde este întâlnit şirul <i>split</i> . De exemplu şirul "as de frt yt uutyh" este scindat prin comanda <code>sir1=strsplit(sir, "t")</code> în şirurile: "as de fr", " y", " uu", "yh". <code>sir1</code> este o listă ce conţine un singur vector <code>sir1[[1]]</code> ce are ca elemente subşirurile rezultat.
<code>paste(..., sep=" ", collapse=NULL)</code>	Converteşte argumentele în şiruri de caractere pe care le concatenează utilizând separatorul <i>sep</i> . Dacă unele argumente sunt vectori se face concatenarea termen cu termen repetând ciclic elementele vectorilor mai scurţi. Dacă De exemplu, <code>paste("a", 1:3, 1:4, sep=" ")</code> are ca rezultat "a11" "a22" "a33" "a14". Dacă <code>collapse=""</code> atunci elementele apar ca un singur şir de caractere. Comanda <code>paste("a", 1:3, 1:4, sep=" ", collapse="")</code> dă ca rezultat: "a11a22a33a14".
<code>toupper(x)</code>	Converteşte la majuscule
<code>tolower(x)</code>	Converteşte la minuscule
<code>LETTERS, letters</code>	Sunt vectori cu caracterele majuscule (minuscul) ale alfabetului latin.

Exemplul 3.6 Urmărind comenzile de mai jos se pot vedea unele variante de aplicare a operaţiilor cu şiruri.

```
> x="avxfdrwyudbcdbdgfjavsjhhdgderwrwy"
> nchar(x)
[1] 32

> x1=strsplit(x,split=NULL)
> x1
[[1]]
[1] "a" "v" "x" "f" "d" "r" "w" "y" "u" "d" "b" "c" "b"
[14] "d" "g" "f" "j" "a" "v" "s" "j" "h" "d" "h" "g" "d"
[27] "e" "r" "w" "r" "w" "y"

> x2=unlist(x1)
> x2
[1] "a" "v" "x" "f" "d" "r" "w" "y" "u" "d" "b" "c" "b"
[14] "d" "g" "f" "j" "a" "v" "s" "j" "h" "d" "h" "g" "d"
[27] "e" "r" "w" "r" "w" "y"

> x3=rev(x2)
> x3
[1] "y" "w" "r" "w" "r" "e" "d" "g" "h" "d" "h" "j" "s"
[14] "v" "a" "j" "f" "g" "d" "b" "c" "b" "d" "u" "y" "w"
[27] "r" "d" "f" "x" "v" "a"

> x4=paste(x3,collapse="")
```

```

> x4
[1] "ywrwredghdhjsvajfgdbcbduywrdfxva"
> x5=substr(x4,3,7)
> x5
[1] "rwred"
> x6=sub("d","D",x)
> x7=gsub("d","D",x)
> x
[1] "avxfdrwyudbcdbdgfjavsjhdhgderwrwy"
> x6
[1] "avxfDrwyudbcdbdgfjavsjhdhgderwrwy"
> x7
[1] "avxfDrwyuDcbDdgfjavsjhDhgDerwrwy"
> x8=unlist(strsplit(x,split="w"))
> x8
[1] "avxfdr"      "yudbcdbdgfjavsjhdhgder"      "r"
[4] "y"
> i=grep("r",x8)
> i
[1] 1 2 3
> x8[i]
[1] "avxfdr"      "yudbcdbdgfjavsjhdhgder"      "r"

```

Constatăm că şirul de comenzi care duc de la x la $x4$ inversează caracterele într-un şir.

3.4. Repartiții de probabilitate

R pune la dispoziția utilizatorilor un mare număr de distribuții de probabilitate. Numele pentru densitatea de probabilitate începe cu d (de exemplu, `dnorm`), pentru funcția de repartiție începe cu p (de exemplu, `pnorm`), pentru cuantile începe cu q (de exemplu, `qnorm`), pentru vectori aleatori începe cu r (de exemplu, `rnorm`).

Densitatea de probabilitate are argumentul opțional `log=FALSE`. Dacă `log=TRUE` atunci se calculează logaritmul din densitate.

Funcția de repartiție are parametrii opționali `lower.tail = TRUE`, `log.p = FALSE`. Dacă `lower.tail=TRUE` atunci funcția de repartiție returnează valoarea $P(X \leq x)$, altfel $P(X > x)$. Dacă `log.p=TRUE` se interpretează că probabilitățile sunt date prin logaritmi lor.

Funcția care calculează cuantilele are parametrii opționali `lower.tail = TRUE`, `log.p = FALSE`. Dacă `lower.tail = TRUE`, atunci funcția aplicată lui q întoarce valoarea minimă x astfel ca $P(X \leq x) \geq q$.

În tabelul de mai jos sunt prezentate câteva distribuții disponibile în R.

Distribuția	Modul de apelare
Normală $\rho(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$	<code>dnorm(x, mean=0, sd=1, log=FALSE)</code> <code>pnorm(x, mean=0, sd=1, lower.tail=TRUE, log.p=FALSE)</code> <code>qnorm(q, mean=0, sd=1, lower.tail=TRUE, log.p=FALSE)</code> <code>rnorm(n, mean, sd)</code> , n este numărul de valori aleatoare

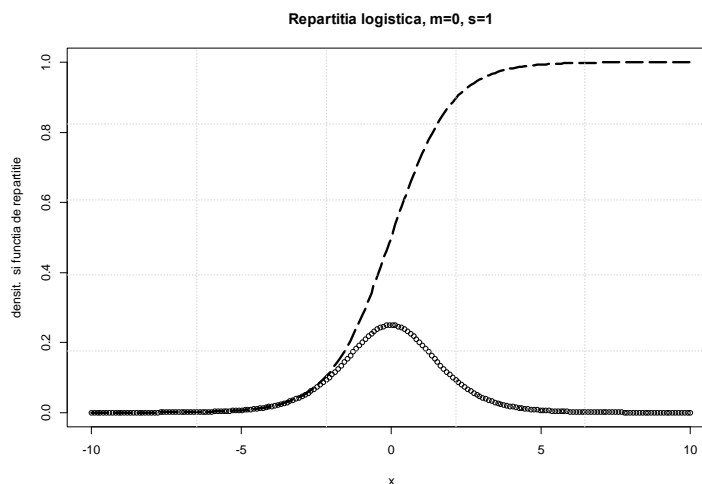
	mean = m, sd=σ
<p>Log-normal</p> $\rho(x) = \frac{1}{\sqrt{2\pi}\sigma x} e^{-\frac{(\log x - \mu)^2}{2\sigma^2}}$	<p>dlnorm(x, meanlog=0, sdlog=1, log=FALSE) plnorm(q, meanlog=0, sdlog=1, lower.tail=TRUE, log.p=FALSE) qlnorm(p, meanlog=0, sdlog=1, lower.tail=TRUE, log.p=FALSE) rlnorm(n, meanlog = 0, sdlog = 1) meanlog = μ, sdlog = σ</p>
<p>Binomială</p> $P(X=x) = C_n^x p^x (1-p)^{n-x}$ <p>$x = 0, 1, 2, \dots, n$</p>	<p>dbinom(x, size, prob, log=FALSE) pbinom(q, size, prob, lower.tail=TRUE, log.p = FALSE) qbinom(p, size, prob, lower.tail=TRUE, log.p = FALSE) rbinom(nr, size, prob), <i>nr</i> este numărul de valori aleatoare size=n, prob=p</p>
<p>Binomială negativă</p> $P(X=x) = \frac{\Gamma(x+n)}{\Gamma(n)x!} p^n (1-p)^x$ <p>$x \in N, n > 0$</p>	<p>dnbinom(x, size, prob, log = FALSE) pnbinom(q, size, prob, lower.tail=TRUE, log.p = FALSE) qnbinom(p, size, prob, lower.tail=TRUE, log.p = FALSE) rnbinom(n, size, prob) size=n, prob=p</p>
<p>Student</p> $\rho(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$	<p>dt(x, df, ncp=0, log=FALSE) pt(q, df, ncp=0, lower.tail = TRUE, log.p = FALSE) qt(p, df, ncp=0, lower.tail = TRUE, log.p = FALSE) rt(n, df, ncp=0), <i>n</i> este numărul de valori aleatoare df = ν, ncp nu este de obicei utilizat</p>
<p>Hi pătrat</p> $\rho(x) = \frac{1}{2^{\nu/2} \Gamma(\nu/2)} x^{\nu/2-1} e^{-x/2},$ <p>$x > 0$</p>	<p>dchisq(x, df, ncp=0, log=FALSE) pchisq(q, df, ncp=0, lower.tail=TRUE, log.p = FALSE) qchisq(p, df, ncp=0, lower.tail=TRUE, log.p = FALSE) rchisq(n, df, ncp=0) df = ν ncp nu se utilizează de obicei, implicit e zero.</p>
<p>Cauchy</p> $\rho(x) = \frac{1}{\pi\gamma \left(1 + \left(\frac{x - x_0}{\gamma}\right)^2\right)}$	<p>dcauchy(x, location=0, scale=1, log=FALSE) pcauchy(q, location = 0, scale = 1, lower.tail = TRUE, log.p = FALSE) qcauchy(p, location = 0, scale = 1, lower.tail = TRUE, log.p = FALSE) rcauchy(n, location = 0, scale = 1) location = x₀, scale = γ</p>
<p>Exponențială</p> $\rho(x) = \lambda e^{-\lambda x}, x > 0, \lambda > 0$	<p>dexp(x, rate = 1, log=FALSE) pexp(q, rate = 1, lower.tail = TRUE, log.p = FALSE) qexp(p, rate = 1, lower.tail = TRUE, log.p = FALSE) rexp(n, rate = 1) rate=λ</p>

<p>Gamma</p> $\rho(x) = \frac{x^{k-1}}{\Gamma(k)a^k} x^{-\frac{x}{a}}, x \geq 0, k > 0, a > 0$	<pre> dgamma(x, shape, scale=1, log=FALSE) pgamma(q, shape, scale = 1, lower.tail = TRUE, log.p = FALSE) qgamma(p, shape, scale = 1, lower.tail = TRUE, log.p = FALSE) rgamma(n, shape, scale = 1) shape=k, scale=a </pre>
<p>Beta</p> $\rho(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$ $a > 0, \beta > 0, x \in (0, 1)$	<pre> dbeta(x, shapel, shape2, ncp = 0, log=FALSE) pbeta(q, shapel, shape2, ncp = 0, lower.tail = TRUE, log.p = FALSE) qbeta(p, shapel, shape2, ncp = 0, lower.tail = TRUE, log.p = FALSE) rbeta(n, shapel, shape2, ncp = 0) shapel=α, shape2=β </pre>
<p>F</p> $\rho(x) = \frac{1}{\beta\left(\frac{a}{2}, \frac{b}{2}\right)} \left(\frac{a}{b}\right)^{a/2} x^{\frac{a}{2}-1} \left(1 + \frac{a}{b}x\right)^{-\frac{a+b}{2}}$ $a > 0, b > 0, x \geq 0, \beta \text{ este funcția beta}$	<pre> df(x, df1, df2, ncp=0, log = FALSE) pf(q, df1, df2, ncp=0, lower.tail=TRUE, log.p = FALSE) qf(p, df1, df2, ncp=0, lower.tail=TRUE, log.p=FALSE) rf(n, df1, df2, ncp=0) df1=a, df2=b </pre>
<p>Geometrică</p> $P(X=x) = p(1-x)^{x-1}, x \in N, p \in [0,1]$	<pre> dgeom(x, prob, log = FALSE) pgeom(q, prob, lower.tail = TRUE, log.p = FALSE) qgeom(p, prob, lower.tail = TRUE, log.p = FALSE) rgeom(n, prob) prob=p </pre>
<p>Hipergeometrică</p> $P(X=x) = \frac{C_m^x C_n^{k-x}}{C_{m+n}^k}$	<pre> dhyper(x, m, n, k, log = FALSE) phyper(q, m, n, k, lower.tail = TRUE, log.p = FALSE) qhyper(p, m, n, k, lower.tail = TRUE, log.p = FALSE) rhyper(nn, m, n, k) </pre>
<p>Poisson</p> $P(X=x) = \frac{\lambda^x e^{-\lambda}}{k!}, \lambda > 0, x \in N$	<pre> dpois(x, lambda, log = FALSE) ppois(q, lambda, lower.tail = TRUE, log.p = FALSE) qpois(p, lambda, lower.tail = TRUE, log.p = FALSE) rpois(n, lambda) lambda = λ </pre>
<p>Uniformă</p> $\rho(x) = \frac{1}{b-a}, x \in [a, b]$	<pre> dunif(x, min=0, max=1, log = FALSE) punif(q, min=0, max=1, lower.tail=TRUE, log.p = FALSE) qunif(p, min=0, max=1, lower.tail=TRUE, log.p = FALSE) runif(n, min=0, max=1) min=a, max=b </pre>

Weibull $\rho(x) = \frac{a}{b} \left(\frac{x}{b}\right)^{a-1} e^{-\left(\frac{x}{b}\right)^a}, x > 0, a > 0, b > 0$	dweibull(x, shape, scale=1, log=FALSE) pweibull(q, shape, scale = 1, lower.tail = TRUE, log.p = FALSE) qweibull(p, shape, scale = 1, lower.tail = TRUE, log.p = FALSE) rweibull(n, shape, scale = 1) shape=a, scale=b
Logistică $\rho(x) = \frac{1}{s} \frac{e^{-\frac{x-m}{s}}}{\left(1 + e^{-\frac{x-m}{s}}\right)^2}$	dlogis(x, location = 0, scale = 1, log = FALSE) plogis(q, location = 0, scale = 1, lower.tail = TRUE, log.p = FALSE) qlogis(p, location = 0, scale = 1, lower.tail = TRUE, log.p = FALSE) rlogis(n, location = 0, scale = 1) location=m, scale=s

Exemplul 3.7. Să realizăm graficul densității și funcția de repartiție pentru distribuția logistică, $m = 0, s = 1$.

```
x=seq(-10,10,length=200)
y=dlogis(x, location = 0, scale = 1, log = FALSE)
y1=plogis(x, location = 0, scale = 1, log = FALSE)
plot(x, y, ylim=range(y,y1), ylab="densitatea si functia
de repartitie", main="Repartitia logistica, m=0,s=1")
lines(x,y1, lwd=3, lty=5), grid(5,5)
```



3.5. Funcții statistice elementare

Funcția	Explicații
<code>mean(x, trim=0, na.rm=FALSE)</code>	Media obiectului x . <code>trim</code> indică ce fracție din valorile dinspre capete se neglijează. Parametrul <code>na.rm</code> indică dacă se elimină NA și NaN din calcul.
<code>sd(x, na.rm=FALSE)</code>	Deviația standard

<code>var(x, y = NULL, na.rm = FALSE, use)</code>	Varianța (dispersia) vectorului, matricei sau <code>data.frame</code> <code>x</code> . Parametrul <code>use</code> este unul din șirurile: "everything", "all.obs", "complete.obs", "na.or.complete" sau "pairwise.complete.obs" și specifică felul în care se tratează valorile lipsă. Vectorul <code>y</code> trebuie să aibă aceeași dimensiune ca <code>x</code> . Implicit <code>y=NULL</code> înseamnă <code>y=x</code> .
<code>median(x, na.rm=FALSE)</code>	Valoarea mediană
<code>cov(x, y = NULL, use = "everything", method = "kendall")</code>	Covarianța vectorilor <code>x</code> și <code>y</code> sau a matricei ori <code>data.frame</code> <code>x</code> (în acest caz se face covarianță între coloane). Parametrul <code>method</code> poate fi "pearson", "kendall", "spearman".
<code>cor(x, y = NULL, use = "everything", method = "kendall")</code>	Coeficientul de corelație al vectorilor <code>x</code> și <code>y</code> sau a matricei ori <code>data.frame</code> <code>x</code> (în acest caz se face corelația între coloane). Parametrul <code>method</code> poate fi "pearson", "kendall", "spearman".
<code>quantile(x, prob, na.rm=FALSE, type=7)</code>	Se determină cuantilele din vectorul de probabilități <code>prob</code> pentru datele numerice din <code>x</code> . Parametrul <code>na.rm</code> controlează dacă se exclud valorile NA și NaN. Parametrul <code>type</code> ia valori întregi între 1 și 9 în funcție de metoda dorită de calcul a cuantilelor.
<code>range(x)</code>	Intervalul [minim, maxim]
<code>sum(x, y), prod(x, y, ...)</code>	Suma elementelor vectorilor <code>x, y, ..</code> respectiv produsul.
<code>diff(x, lag=1)</code>	Diferențe finite $x_{i+lag} - x$, cu pasul <code>lag</code> . Implicit este <code>lag=1</code> , iar rezultatul are lungimea cu 1 mai mică. Dacă <code>x</code> este matrice se fac diferențele pe componente.
<code>min(x)</code>	Minimul
<code>max(x)</code>	Maximul
<code>summary(obiect, ...)</code>	Este o funcție generică al cărei rezultat depinde de obiectul căreia i se aplică. Pentru un vector numeric se dau minimul, maximum, media, mediana, prima și a treia cvartilă. Pentru o matrice sau <code>data frame</code> se face același lucru pentru coloane.
<code>scale(x, center=TRUE, scale=TRUE)</code>	Se înlocuiește x_i cu $\frac{x_i - center}{scale}$. Dacă <code>center=TRUE</code> atunci se ia <code>center</code> egal cu <code>mean(x)</code> , iar dacă <code>scale=TRUE</code> se ia <code>scale=sd(x)</code> . Dacă <code>x</code> este o matrice se aplică procedura pentru fiecare coloană. <code>Center</code> și <code>scale</code> trebuie să fie în acest caz un vector de centre, respectiv un vector de parametri de scalare.

Exemplul 3.8. Se dă un vector numeric x și se cere un sumar al caracteristicilor statistice elementare.

```
> x<-c(1,2,-1,3,4,10,6,4,-2,6,8,3,5)
> summary(x)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-2.000   2.000   4.000   3.769   6.000  10.000
```

3.6. Funcții de mapare

Prin aceste funcții se evită folosirea unor instrucțiuni de ciclare. O anumită funcție este aplicată repetat unor date extrase din argumentele instrucțiunii.

Funcția	Explicații
<code>apply(x, MARGIN, FUN, ...)</code>	x este un array, iar MARGIN un vector cu pozițiile pe care sunt indicii care se ciclează, FUN este funcția care se aplică acelui array cu indicii din MARGIN fixați, ... sunt parametrii adiționali pentru FUN. De exemplu, dacă x este un array cu mai mult de două dimensiuni, să zicem trei, comanda <code>y=apply(x, c(1,2), sum)</code> ne dă o matrice unde <code>y[i,j]=sum(x[i,j,])</code> , pentru orice i,j între limitele permise, iar comanda <code>y=apply(x,c(1,3),sum)</code> ne dă <code>y[i,j]=sum(x[i,j,])</code> .
<code>lapply(x, FUN, ...)</code>	Aplică funcția FUN fiecărui element al unui vector sau listă.
<code>sapply(x, FUN, ..., simplify = TRUE, ...)</code>	La fel ca <code>lapply</code> doar că rezultatul se convertește (simplifică) la un array.
<code>replicate(n, expr, simplify = "array")</code>	Evaluarea de n ori a expresiei <code>expr</code> . Dacă x este un vector atunci <code>y=replicate(3,mean(x))</code> ne dă un vector cu trei componente egale cu media lui x .
<code>tapply(x, INDEX, FUN = NULL, ..., simplify = TRUE)</code>	Se aplică funcția FUN elementelor din vectorul x grupate după lista de factori INDEX. Dacă <code>simplify=TRUE</code> se încearcă convertirea la array a rezultatului.
<code>mapply(FUN, ..., MoreArgs = NULL, simplify = TRUE, use.names = TRUE)</code>	Se aplică funcția FUN, funcție de mai multe argumente, primelor elemente ale argumentelor..., pe urmă elementelor de pe locurile doi ale argumentelor, etc. ... poate fi o înșiruire de vectori da aceeași lungime, data frame cu același număr de coloane sau liste cu același număr de componente. De exemplu dacă <code>z1=c(1,2,3,4)</code> , <code>z2=c(3,4,5,6)</code> atunci <code>mapply(sum,z1,z2)</code> ne dă 4 6 8 10.

Exemplul 3.9. Să aplicăm funcțiile `apply`, `lapply`, `sapply` și `mapply` unor matrice sau vectori. Rezultatele se pot vedea după fiecare comandă.

```
> x=matrix(1:12,ncol=4)
> x
     [,1] [,2] [,3] [,4]
[1,]    1    4    7   10
[2,]    2    5    8   11
[3,]    3    6    9   12
```

```

> x1=apply(x,MARGIN=1,FUN=sum)
> x1
[1] 22 26 30
> x2=apply(x,MARGIN=2,FUN=sum)
> x2
[1] 6 15 24 33
> x3=lapply(x,FUN=sqrt)
> x3
[[1]]
[1] 1

[[2]]
[1] 1.414214

[[3]]
[1] 1.732051

[[4]]
[1] 2

[[5]]
[1] 2.236068

[[6]]
[1] 2.44949

[[7]]
[1] 2.645751

[[8]]
[1] 2.828427

[[9]]
[1] 3

[[10]]
[1] 3.162278

[[11]]
[1] 3.316625

[[12]]
[1] 3.464102

> x4=sapply(x,FUN=sqrt)
> x4
[1] 1.000000 1.414214 1.732051 2.000000 2.236068
[6] 2.449490 2.645751 2.828427 3.000000 3.162278
[12] 3.316625 3.464102

> y=matrix(seq(from=0,by=10,length=12),ncol=4);
> y
      [,1] [,2] [,3] [,4]
[1,]    0   30   60   90
[2,]   10   40   70  100
[3,]   20   50   80  110

> z=mapply(sum,x,y)
> z

```

```

[1] 1 12 23 34 45 56 67 78 89 100 111 122
> z1=matrix(z,ncol=ncol(x))
> z1
      [,1] [,2] [,3] [,4]
[1,]    1   34   67  100
[2,]   12   45   78  111
[3,]   23   56   89  122

```

Se constată că funcția `lapplay` are ca rezultat o listă. Pentru a obține un vector trebuie utilizată funcția `sapply` în loc de `lapplay` sau aplicată funcția `unlist` rezultatului întors de `lapplay`.

Exemplul 3.10. Să calculăm produsul de convoluție p al vectorilor x și y . Prin

definiție $p_k = \sum_{i+j=k} x_i y_j$. Indicii vectorului x apar în vectorul i_x , ai vectorului y în

i_y , iar ai vectorului p în i_p .

```

> x=1:5;y=c(1,2,2);ix=-2:2; iy=0:2
> x
[1] 1 2 3 4 5
> ix
[1] -2 -1 0 1 2
> y
[1] 1 2 2
> iy
[1] 0 1 2

> r=as.vector(outer(x,y,"*"))
> ir=as.vector(outer(ix,iy,"+"))
> r
[1] 1 2 3 4 5 2 4 6 8 10 2 4 6 8 10
> ir
[1] -2 -1 0 1 2 -1 0 1 2 3 0 1 2 3 4

> fir=factor(ir)
> p=tapply(r,fir,sum)
> ip=as.numeric(unlist(dimnames(p)))
> p
-2 -1 0 1 2 3 4
 1  4  9 14 19 18 10
> ip
[1] -2 -1 0 1 2 3 4

> dimnames(p)<-NULL
> p
[1] 1 4 9 14 19 18 10

```

3.7. Funcții pentru citire și scriere

Descrierile care urmează sunt incomplete, dar suficiente în aplicații uzuale. Prin apel la `help` se pot obține informații suplimentare.

Funcția	Explicații
<code>print</code>	<code>print(x, ...)</code> Se tipăresc la consolă datele <code>x</code> . Printre opțiuni se numără <code>digits=</code> numărul maxim de cifre semnificative ce se tipăresc, <code>justify="left", "right", "centre", "none"</code> .
<code>cat</code>	<code>cat(..., file="", sep="", fill=FALSE, append=FALSE)</code> Vectorii și numele din ... sunt scrise în fișierul <code>file</code> (dacă nu apare în instrucțiune atunci se scrie la consolă) utilizând separatorul <code>sep</code> . <code>fill</code> este un număr pozitiv care dă lungimea unei linii în fișier sau <code>FALSE</code> când se crează linie nouă doar la apariția caracterului “\n”. Dacă <code>append=TRUE</code> se adaugă datele la fișierul existent.
<code>read.csv</code> , <code>read.table</code> , <code>read.csv2</code> , <code>read.delim</code> , <code>read.delim2</code>	Prin aceste instrucțiuni se citesc date de tip <code>data.frame</code> din fișiere text. Pe fiecare linie din fișier trebuie să fie o linie din tabel. Dacă <code>header=TRUE</code> atunci prima linie conține numele coloanelor. Parametrul <code>row.names</code> este fie un număr ce indică pe ce coloană este numele liniilor fie un sir de caractere cu numele coloanei. <code>read.csv(file, header=TRUE, sep=" ", quote="\\"", dec=".")</code> <code>read.table(file, header=FALSE, sep=" ", quote="\'", dec=".", row.names, col.names, , nrows=-1, skip=0, blank.lines.skip=TRUE, comment.char="#",)</code> Parametrii din instrucțiuni se explică singuri prin nume (<code>dec</code> este semnul pentru separarea zecimalelor). Dacă lipsesc se pun valorile implicite. <code>read.csv2</code> , <code>read.delim</code> , <code>read.delim2</code> sunt la fel ca <code>read.csv</code> cu alți separatori între câmpuri și alți delimitatori pentru zecimale.
<code>readLines</code>	Se utilizează pentru a citi linii de la o conexiune (de exemplu dintr-un fișier)
<code>write</code> , <code>write.table</code> , <code>write.csv</code> , <code>write.csv2</code>	<code>write(x, file, ncolumns, append=FALSE, sep=" ")</code> <code>write.table(x, file="", append=FALSE, quote=TRUE, sep=" ", eol="\n", na="NA", dec=".", row.names=TRUE, col.names = TRUE)</code> Parametrii care nu se pun efectiv în instrucțiuni se iau implicit. Cu instrucțiuni se scriu date de tip <code>data.frame</code> în fișiere text.
<code>scan</code>	<code>scan(file="", what, nmax=-1, n=-1, sep=" ", skip=0, nlines=0)</code> , unde <code>file</code> este fișierul, <code>what</code> este tipul de date din fișier: <code>logical</code> , <code>integer</code> , <code>numeric</code> , <code>complex</code> , <code>character</code> , <code>raw</code> , <code>list</code> (implicit <code>numeric</code>), <code>sep</code> este separatorul de date din fișier (implicit <code>blanc</code>), <code>skip</code> este numărul de linii de la începutul fișierul de ignorat, <code>nlines</code> este numărul max.de linii ce se citesc.

format	format(x, digits, nsmall, justify, width) Se specifică felul în care se va tipări x, digits este numărul de cifre, nsmall este numărul de cifre după virgulă (≤ 20), justify="left","right","centre","none", width este lățimea min.
save	save(...,list,file,ascii=FALSE), unde ... este lista de variabile ce se salvează, list este la fel liste de variabile ce se salvează, file este fișierul în care se salvează, ascii este TRUE sau FALSE.
save.image	save.image(file=".RData", ascii=FALSE, compress, safe = TRUE) Compress poate fi "gzip", "bzip2" or "xz" și specifică tipul de compresie al datelor. Prin această comandă se salvează toate datele din spațiul de lucru.
load	load(file), unde file este fișierul de unde se restaurează datele salvate cu save sau save.image
ls()	Rezultatul este o listă cu variabilele active în programul curent
rm(list)	Prin această funcție sunt șterse din memorie variabilele din lista dată de parametrul list. Instrucțiunea rm(list=ls()) șterge toate variabilele din memorie.

În capitolul despre tipuri de date s-au dat mai multe exemple de citire sau scriere a datelor. În exemplul următor ilustrăm utilizarea funcțiilor paste, cat, print.

Exemplul 3.11 Un vector este scris la consolă cu componentele numite detaliat.

```
> x=1:10
> y=sqrt(x)
> et=paste("y",x,sep="")
> et
[1] "y1" "y2" "y3" "y4" "y5" "y6" "y7" "y8" "y9" "y10"
> cat(paste(et,"=",format(y,digits=3),";",sep=""),"\n")
y1=1.00; y2=1.41; y3=1.73; y4=2.00; y5=2.24; y6=2.45;
y7=2.65; y8=2.83; y9=3.00; y10=3.16;
> print(y, digits=3)
[1] 1.00 1.41 1.73 2.00 2.24 2.45 2.65 2.83 3.00 3.16
```

3.8. Funcții de creare de date

Printre funcțiile de creare de date discutate a capitolul despre tipuri de date amintim: rep, seq, c, matrix, array, table, cbind, rbind. Aceste funcții au fost utilizate în mai multe rânduri în capitolul despre tipurile de date.

Funcția	Exemple
seq(from, to, by, length.out, along.with=NULL, ...)	X=seq(0,1,length=100) Y=seq(from=1, to=4, by=0.1) Z=seq(from=10,by=0.3,along=X)
rep(x, ...). La opțiuni avem times, length.out, each (la each se repetă fiecare element al lui x)	rep(x,10); rep(x,times=10); rep(1:3,length=5); rep(1:5,each=2)

	<code>x<-c(2,3,4)</code>
<code>cbind, rbind</code> combina pe coloanele (liniile) unor vector, ale unei matrice sau data.frame	<code>x=1:3; y=4:6; z<-cbind(x,y);</code> <code>u=rbind(x,y)</code>
<code>matrix, array</code> permit crearea de matrice sau blocuri multidimensionale	<code>x=matrix(1:10, ncol=2)</code> <code>y=array(1:24, dim=c(2,3,4))</code>

Exemplul 3.12. Ilustrăm cum apar rezultatele la unele modele de funcții din tabelul de mai sus.

```
> rep(1:3,length=5)
[1] 1 2 3 1 2
> rep(1:5,each=2)
[1] 1 1 2 2 3 3 4 4 5 5
> x=1:3; y=4:6; z<-cbind(x,y);
> z
      x y
[1,] 1 4
[2,] 2 5
[3,] 3 6
> t=array(1:24, dim=c(2,3,4))
> t
, , 1

      [,1] [,2] [,3]
[1,]     1     3     5
[2,]     2     4     6

, , 2

      [,1] [,2] [,3]
[1,]     7     9    11
[2,]     8    10    12

, , 3

      [,1] [,2] [,3]
[1,]    13    15    17
[2,]    14    16    18

, , 4

      [,1] [,2] [,3]
[1,]    19    21    23
[2,]    20    22    24
```

3.9. Funcții de conversie

Conversia datelor dintr-un tip în altul se poate face prin comenzi ca `as.numeric`, `as.integer`, `as.character`, `as.factor`, `as.function`, `as.vector`, `cut` (convertește un vector numeric într-un factor), `strtoi` (convertește șiruri de caractere în întregi), `toString` (convertește un obiect R într-un șir de caractere), `data.matrix` (convertește un obiect `data.frame` într-o matrice) etc. Dacă nu e posibilă conversia atunci rezultatul rămâne nemodificat.

Exemplul 3.13. Din doi vectori de aceeași dimensiune se construiește o `data.frame` care se încearcă a se converti în mai multe feluri dintre care unele sunt imposibile.

```
> x=1:5
> y=letters[1:5]
> a=data.frame(x,y)
> a
  x y
1 1 a
2 2 b
3 3 c
4 4 d
5 5 e
> class(a)
[1] "data.frame"
> b=as.character(a)
> b
[1] "1:5" "1:5"
> c=as.matrix(a)
> c
      x y
[1,] "1" "a"
[2,] "2" "b"
[3,] "3" "c"
[4,] "4" "d"
[5,] "5" "e"
> as.numeric(c)
[1] 1 2 3 4 5 NA NA NA NA NA
Warning message:
NAs introduced by coercion
```

3.10. Funcții pentru teste statistice

Pentru teste statistice se poate apela în R la funcții specializate ca: *t.test*, *binom.test*, *prop.test*, *var.test*, *wilcox.test*, *shapiro.test*, *bartlett.test*, *Box.test*, *chisq.test*, *fisher.test*, *friedman.test*, *ks.test*, etc. Informații despre aceste funcții, parametrii de intrare și rezultatele întoarse, se pot găsi în paginile de ajutor din R. Valoarea de bază întoarsă de un test este `p.value`, notată în exterior ca *p-value*. Dacă este mai mare decât pragul de risc definit 1-nivelul de încredere (în teste nivelul de încredere este numit `conf.level`) atunci ipoteza este acceptabilă cu o eroare de tipul I sub pragul de risc.

Exemplul 3.14. Din 1000 aruncări cu o monedă, fața a apărut de 490 ori. Să se testeze dacă probabilitatea $p=1/2$ de apariție a feței este admisibilă cu un prag de risc de 0.05. Următoarea comandă în R face acest test.

```
> rez=binom.test(490, 1000, p=0.5, conf.level=0.95)
> rez
      Exact binomial test
data:  490 and 1000
number of successes = 490, number of trials = 1000,
p-value = 0.548
alternative hypothesis: true probability of success is not
equal to 0.5
95 percent confidence interval:
 0.4585849 0.5214742
sample estimates:
probability of success
              0.49
```

În cazul nostru $p\text{-value}=0.548 > 0.05=\text{pragul de risc}=1-\text{conf.level}$, deci ipoteza că $p=0.5$ este admisibilă. În plus, avem și o estimare a lui p sub forma unui interval de încredere cu încrederea $\text{conf.level}=0.95$.

Exemplul 3.15. Utilizând testul Kolmogorov-Smirnov să se decidă dacă datele dintr-un vector sunt repartizate conform unei distribuții date (în cazul de față $N(1,2)$).

```
> x <- runif(500,-1,1)
> rez=ks.test(x, "pnorm",1,2)
> rez
      One-sample Kolmogorov-Smirnov test
data:  x
D = 0.5005, p-value < 2.2e-16
alternative hypothesis: two-sided
```

Am generat un vector aleator uniform cu valori în $[-1,1]$. Testul `ks` a detectat că datele nu provin dintr-o distribuție normală cu media 1 și abaterea standard 2. Acest lucru se reflectă prin faptul că $p\text{-value}$ este foarte mic, mult mai mic decât pragul de risc uzual 0.05. Dacă schimbăm datele de testare cu un vector repartizat $N(1,2)$, testul `ks` arată că ipoteza este acceptabilă ($p\text{-value}$ este $0.7422 > 0.05=\text{pragul uzual de risc}$)

```
> x <- rnorm(500,1,2)
> rez=ks.test(x, "pnorm",1,2)
> rez
      One-sample Kolmogorov-Smirnov test
data:  x
D = 0.0304, p-value = 0.7442
alternative hypothesis: two-sided
```

3.11. Funcții pentru determinarea duratei de execuție a unor porțiuni de program

Funcția	Exemple
<code>system.time(expr)</code>	<p>Este determinată durata de execuție a expresiei <code>expr</code>. De exemplu</p> <pre>> system.time(for(i in 1:10000) x <- mean(rnorm(1000, mean=4, sd=5))) user system elapsed 1.72 0.00 1.72</pre>
<code>proc.time()</code>	<p>Este determinat timpul cât calculatorul cheltuie cu anumite proceduri. Se utilizează ca de exemplu:</p> <pre>ptm <- proc.time() a=matrix(runif(1000000),ncol=1000) b=runif(1000) x=solve(a,b) proc.time() - ptm user system elapsed 0.69 0.08 1.36</pre>

Numărul funcțiilor predefinite în R este mult mai mare decât am descris noi aici și pot apărea funcții noi în versiunile ulterioare. Funcțiile descrise în tabelele de mai sus fac parte din pachetul *base* și din pachetul *stats*, și foarte probabil vor rămâne și în versiunile ulterioare de R. Pachetele adiționale conțin multe funcții suplimentare care sunt orientate spre rezolvarea unor probleme specifice. Vom mai descrie ulterior unele funcții din pachetul *tseries*, funcții orientate spre lucrul cu serii temporale.

Capitolul 4

Funcții grafice în R

În R se pot executa grafice de calitate pentru datele admise. Aceste grafice sunt trimise pe ecran sau în fișiere de diverse tipuri. Pentru a realiza acest lucru trebuie deschis un “device”. Comanda *windows()* deschide ecranul ca dispozitiv grafic, comanda *jpeg()* deschide ca dispozitiv grafic un fișier în care se va scrie informația grafică în format jpeg, *postscript()* deschide un fișier de tip postscript, *pdf()*, *png()*, *tiff()*, *bitmap()*, *win.metafile()*, etc. Aceste comenzi au ca parametru *file* numele fișierului unde se salvează graficul, *width*=lățimea în pixeli, *height*=înălțimea în pixeli. Numai un dispozitiv este activ la un moment, trecerea de la unul la altul se face *dev.next()*, *dev.prev()*, *dev.set(which=k)*. Listarea dispozitivelor grafice se face cu comanda *dev.list()*, închiderea unui dispozitiv cu *dev.off(k)*, iar închiderea tuturor cu *graphics.off()*. Când s-a închis un dispozitiv diferit de ecran, imaginea din el este salvată în fișierul specificat în parametrul *file*. După comanda *graphics.off()* trebuie deschis un nou dispozitiv pentru a putea tipări grafic, de regulă *windows()* pentru ecran. Un clic dreapta pe fereastra cu reprezentarea grafică de pe ecran deschide un meniu cu opțiuni de salvare sau copiere a conținutului în clipboard.

Rutinele grafice sunt cuprinse în modulul de bază: “base” precum și în două pachete suplimentare “grid” și “lattice”. Rutinele de bază se împart în trei categorii: rutine de nivel înalt “high level” în urma cărora se crează o fereastră grafică nouă în care se pun graficele specificate prin argumentele rutinelor; rutine de nivel jos “low level” prin care la graficul existent se adaugă informații suplimentare ca text, puncte, linii; grafice interactive prin care utilizatorul poate interacționa cu graficul prin intermediul mouse-ului pentru a obține informații despre valorile din diverse poziții de pe grafic sau pentru a adăuga informații în diverse poziții de pe grafic.

Rutinele de nivel înalt din modulul de bază sunt:

Instrucțiunea	Explicații
plot	<i>plot(x,y,...)</i> , <i>plot(xy,...)</i> , tipărește punctele de coordonate ($x[i], y[i]$) utilizând opțiunile grafice ... (vezi mai jos). <i>plot(f,x)</i> unde <i>f</i> este un factor pentru <i>x</i> tipărește un <i>boxplot</i> pentru <i>x</i> grupat după factorul <i>f</i> . <i>plot(x~y)</i> face graficele coloanelor $y^{<i> </i>}$ în funcție de <i>x</i> . <i>plot(x)</i> pentru o variabilă care are metoda <i>plot</i> produce un grafic specific acelui tip de variabilă.
pairs	<i>pairs(x)</i> unde <i>x</i> este o matrice face toate graficele $(x^{<i> </i>}, x^{<j> </j>})$ unde $x^{<i> </i>}$ și $x^{<j> </j>}$ sunt coloane ale lui <i>x</i>
coplot	<i>coplot(x~y,f)</i> execută graficele lui <i>y</i> depinzând de <i>x</i> , pentru fiecare nivel al factorului <i>f</i>
hist	<i>hist(x,...)</i> tipărește histograma frecvențelor valorilor din vectorul <i>x</i> . Opțiunea frecvent utilizată este <i>nclass=n</i> pentru a specifica în câte clase se împart datele din <i>x</i> . Cu opțiunea <i>probability=TRUE</i> histograma va reprezenta probabilitățile nu frecvențele.

boxplot	<i>boxplot(x,...)</i> sau <i>boxplot(x~y,data=z,...)</i> . În primul caz se execută un <i>boxplot</i> pentru valorile din vectorul sau lista de vectori <i>x</i> , iar în al doilea caz se execută un <i>box plot</i> pentru valorile din coloana <i>x</i> grupate după factorul <i>y</i> , date ce se află în variabila <i>z</i> ce trebuie să aibă date numite <i>x</i> și <i>y</i> (de exemplu <i>z</i> este <i>data.frame</i> iar <i>x</i> , <i>y</i> sunt două coloane în <i>z</i>).
qqnorm, qqplot, qqline	<i>qqnorm(x)</i> reprezintă grafic cuantilele lui <i>x</i> față de cele ale distribuției normale standard, <i>qqplot(x,y)</i> reprezintă grafic cuantilele lui <i>x</i> față de cele ale lui <i>y</i> , <i>qqline(x)</i> reprezintă pe lângă graficul ca în <i>qqnorm(x)</i> și linia ce aproximează graficul.
barplot	<i>barplot(x,...)</i> , unde ... sunt opțiuni, tipărește sub formă de bare verticale valorile lui <i>x</i>
dotchart	<i>dotchart(x,...)</i> la fel ca <i>barplot(x,...)</i> face graficul valorilor lui <i>x</i> sub formă de puncte
image	<i>image(x,y,z,...)</i> , <i>image(z,...)</i> reprezintă matricea <i>z</i> prin culori; <i>x[i]</i> și <i>y[j]</i> reprezintă coordonatele centrului dreptunghiului unde se reprezintă prin culoare valoarea <i>z[i,j]</i>
contour, filled.contour	<i>contour(x,y,z, nlevels=n,...)</i> reprezintă grafic liniile de contur ale valorilor <i>z[i,j]</i> din matricea <i>z</i> , având pe axe gradațiile date de vectorii <i>x</i> , respectiv <i>y</i> . Opțiuni suplimentare se pot afla cu <i>help(contour)</i> . Analog pentru <i>filled.contour</i> , unde spațiul între două linii de contur e colorat.
persp	<i>persp(z)</i> , <i>persp(x,y,z)</i> , <i>persp(x,y,z,theta=a,phi=b,...)</i> desenează o suprafață pe care sunt punctele de coordonate (<i>x[i]</i> , <i>y[j]</i> , <i>z[i,j]</i>). Direcția de vedere a suprafeței se ajustează prin parametrii <i>theta</i> și <i>phi</i> (în grade), iar <i>ltheta</i> și <i>lphi</i> dau direcția de unde vine lumina pe suprafață, <i>r</i> dă distanța de la care se vede (<i>r</i> mare e ca vederea de la infinit). Culoarea este controlată prin parametrul opțional <i>col</i> . Pentru o variație continuă a culorii <i>col</i> trebuie să fie un vector de culori de dimensiune egală cu numărul de fațete (vezi exemplul de mai jos). Pentru a obține o senzație de relief mai puternică se utilizează opțiunea <i>shade=număr</i> de obicei între 0 și 1, <i>thicktype</i> poate fi "simple" sau "detailed", <i>axes</i> este TRUE sau FALSE după cum se pun sau nu marcasele pe axe, <i>xlim</i> , <i>ylim</i> , <i>zlim</i> dau valorile limită pe axe, valori care sunt cuprinse în grafic, <i>main</i> și <i>sub</i> sunt titlul și subtitlul, <i>xlab</i> , <i>ylab</i> , <i>zlab</i> sunt titlurile axelor. Dacă <i>scale</i> = TRUE atunci nu se păstrează proporția pe cele trei axe.

Prin tipărirea cu rutine de nivel jos se adaugă elemente noi la graficul existent. Mai jos sunt descrise unele rutine grafice și unii parametri ai acestora. Restul se pot afla prin comanda *help*.

Rutina grafică	Explicații
points	<i>points(x,y,...)</i> tipărește puncte în coordonatele specificate de <i>x</i> și <i>y</i>
lines	<i>lines(x,y)</i> adaugă linii între punctele de coordonate specificate în vectorii <i>x</i> , <i>y</i> .
text	<i>text(x,y,text,...)</i> determină scrierea începând din poziția (<i>x[i]</i> , <i>y[i]</i>) a caracterelor din <i>text[i]</i> . Un parametru auxiliar numit <i>pos</i> ce poate avea valorile 1, 2, 3, 4 controlează dacă textul este sub, la stânga, deasupra sau la dreapta punctului (<i>x</i> , <i>y</i>).

mtext	<i>mtext(text,...)</i> scrie un text pe o margine a graficului. Printre parametrii opționali : <i>side</i> pentru a controla marginea (1=jos, 2=stânga, 3=sus, 4=dreapta), <i>col</i> este culoarea pentru text, <i>line</i> =întreg controlează depărtarea față de marginea corespunzătoare a graficului (în linii).
abline	<i>abline(a=valoare1, b=valoare2)</i> trasează graficul lui $y=a+bx$, <i>abline(h=valoare)</i> trasează linia orizontală prin $y=h$, <i>abline(v=valoare)</i> trasează linia verticală prin $x=v$., <i>abline(coef=vector)</i> trasează linia care are în <i>coef</i> ordonata la origine și panta, <i>abline(obiect)</i> trasează o linie de parametrii <i>coef</i> din obiectul <i>obiect</i> (dacă există).
axis	<i>axis(side,...)</i> trasează o axă în funcție de parametrul <i>side</i> (1=jos, 2=stânga, 3=sus și 4=dreapta). Există mai mulți parametri opționali.
segments	<i>segments(xo,yo,xd,yd,...)</i> trasează segmente de la (xo,yo) la (xd,yd). Dacă xo, yo, xd, yd sunt vectori atunci se trasează mai multe segmente. Printre parametrii opționali amintim <i>col</i> care este egal cu indicele culorii segmetului (<i>col</i> poate fi un nume de culoare).
arrows	<i>arrows(xo,yo,xd,yd,...)</i> este ca și rutina <i>segments</i> de mai sus doar că trasează săgeți de la (xo,yo) la (xd,yd)
rect	<i>rect(xleft, ybottom, xright, ytop,...)</i> traează o mulțime de dreptunghiuri în care <i>xleft</i> reprezintă coordonatele lor x stânga, <i>ybottom</i> reprezintă coordonatele lor y jos, etc. Printre parametrii opționali menționăm <i>col</i> un vector cu indicii de culoare din paleta curentă, <i>density</i> un vector cu numărul de linii de hașurare, <i>lty</i> un vector cu tipurile de linie de hașură, <i>angle</i> un vector cu unghiurile față de Ox ale liniilor de hșurare.
polygon	<i>polygon(x,y,...)</i> trasează un poligon cu vârfurile consecutive date de vectorii x și y. Printre opțiuni sunt importante <i>col</i> =culoarea cu care se umple interiorul poligonului, <i>density</i> =numărul de linii cu care se hașurează poligonul, <i>angle</i> =unghiul cu axa Ox al liniilor de hașurare (în grade)
box	<i>box(lty=tip, col=culoare)</i> determină trasarea unui dreptunghi în jurul graficului cu tipul de linie <i>tip</i> și culoarea <i>culoare</i> .
grid	<i>grid(nx,ny, col=culoare, lty=tiplinie)</i> datermină trasarea unui caroi aj cu culoarea <i>culoare</i> și tipul de linie <i>tiplinie</i> . Se mai poate introduce parametrul <i>lwd</i> care dă grosimea liniilor de caroi aj. Pe axa Ox sunt <i>nx</i> diviziuni, iar pe axa Oy sunt <i>ny</i> diviziuni.
legend	<i>legend(x,y, legenda, col=vector, text.col=vector, lty=vector, pch=vector, bg=,...)</i> unde (x,y) este poziția pe grafic al colțului stânga sus al legendei (se poate indica poziția prin cuvinte "topleft" sau "bottomright" sau prin instrucțiunea <i>locator()</i> caz în care este aleasă cu mausul), <i>legenda</i> este un vector cu elemente text (se poate da textul din legendă prin <i>legend</i> =textul , <i>col</i> este un vector cu indicii culorilor din paleta curentă de culoare pentru culorile marcajelor din legendă, <i>text.col</i> este un vector cu indicii de culoare pentru textele din legendă, <i>lty</i> este un vector de întregi cu indicii tipurilor de linie din grafice, <i>pch</i> este un vector cu indicii tipurilor de linie din legendă, <i>bg</i> este indicele culorii de fundal din legendă. Culorile, tipurile de linie sau tipurile de puncte pot fi specificate și prin cuvinte)

title	<i>title(main=text1,sub=text2,xlab=text3,ylab=text4,...)</i> determină scrierea unui titlu(main), subtitlu(sub), etichete pentru axe(xlab, ylab) care sunt luate din text1, text2, etc.
-------	---

Pentru diverși parametri grafici R menține o listă numită *par* cu valorile acestora. Cu comanda *par()* putem afla valorile acestor parametri. Cei mai utilizați parametri grafici sunt în tabelul următor:

Parametrul grafic	Valoare implicită	Explicații
xlog	FALSE	Dacă este TRUE pe axa x se utilizează scara logaritmică
ylog	FALSE	Analog cu xlog dar pentru axa y
pch	1	<i>pch</i> este un vector valori întregi care specifică tipul de punct de pe grafic. Valorile întregi sunt de regulă între 0 și 255. După epuizarea tipurilor din <i>pch</i> ele se se repetă.
col	"black"	<i>col</i> este un vector de întregi care dă indexul culorii curente din paleta curentă de culori. Culoarea se poate exprima și prin cuvinte (există 667 de culori care au nume care se pot afla cu comanda <i>colors()</i>). Culorile se repetă ciclic atunci când se reprezintă diverse obiecte grafice.
col.axis	"black"	Culoarea axelor
col.lab	"black"	Culoarea pentru etichete pe axe
col.main	"black"	Culoarea pentru titlu
col.sub	"black"	Culoarea pentru subtitlu
cex, cex.axis, cex.lab, cex.main, cex.sub	1	Un factor de amplificare pentru mărimea textului față de mărimea implicită. Doar <i>cex.main</i> are valoarea implicită 1.2
adj	-1	Se controlează cum apare textul justificat față de punctul de referință: 0 aliniat la stânga, 1 aliniat la dreapta, 0.5 centrat orizontal
font, font.axis, font.lab, font.main, font.sub	1	Un întreg ce reprezintă fontul pentru text, axe, etichete, titlu, subtitlu.
fg, bg	"black", "white"	Culorile pentru foreground respectiv background
lty	"solid"	Tipul de linie: 0=blank, 1=solid (default), 2=dashed, 3=dotted, 4=dotdash, 5=longdash, 6=twodash
lwd	1	Grosimea liniei. Depinde de tipul de dispozitiv grafic pe care se tipărește
mai	1.02 0.82 0.82 0.42	Un vector ce conține dimensiunile în inci ale marginilor în ordinea: jos, stânga, sus, dreapta
mar	5.1 4.1 4.1 2.1	La fel ca la <i>mai</i> dar dimensiunea este în linii
mex	1	Factor de amplificare pentru mărimea fontului cu

		care se scrie pe margine
mfcoll	1 1	Vector de tipul c(nr,nc). Ecranul va fi împărțit în $nr \times nc$ regiuni. Figurile consecutive vor fi scrise pe coloană.
mfrow	1 1	Vector de tipul c(nr,nc). Ecranul va fi împărțit în $nr \times nc$ regiuni. Figurile consecutive vor fi scrise pe linie.
mfg	1 1 1 1	Un vector de tipul c(i,j,nr,nc) prin care se specifică în ce regiune de pe ecran se scrie. Se acceptă și valori de tipul c(i,j)
fig		Printr-un vector de patru numere între 0 și 1 se controlează poziția figurii în pagină față de colțul stânga-jos.
xaxs, yaxs		Controlează stilul marcajelor pe axe și pot fi "r", "i", "e", "s", "d".
xaxp, yaxp	0 1 5	Au forma c(n1,n2,n) unde n1 și n2 reprezintă coordonatele extreme ale punctelor de marcaj (tick marks) de pe axe, iar n este numărul dorit de marcaje
lab	5 5 7	Vector de tipul c(n1,n2,n): n1, n2 reprezintă numerele dorite de marcaje pe axe, iar nu e implementat.
las		Este un întreg 0, 1, 2, sau 3, și controlează felul în care se trasează marcajele pe axe.
xaxt, yaxt	"s"	Controlează tipul axelor "n", "l", "t", "s".
usr	0 1 0 1	Vector de tipul c(x1,x2,y1,y2) pentru valorile extreme pe coordonate
tcl	-0.5	Lungimea marcajelor pe axe ca fracție din lățimea unui rând de text.
srt	0	Unghiul de rotație al șirurilor de caractere

Unii parametri sunt admiși ca date opționale de unele rutine grafice. La apelul rutinei respective, parametrul considerat se introduce sub forma *rutina(..., parametru=valoare,...)*. După execuția comenzii grafice valoarea parametrului revine la cea dinainte de instrucțiunea grafică. Dacă vrem ca modificarea să fie permanentă (până la o nouă modificare) atunci vom utiliza comanda *par(parametru=valoare,...)*. Unii parametri pot fi modificați doar prin instrucțiunea par. Aceștia sunt: "ask", "fig", "fin", "lheight", "mai", "mar", "mex", "mfcoll", "mfrow", "mfg", "new", "oma", "omd", "omi", "pin", "plt", "ps", "pty", "usr", "xlog", "ylog", "ylbias".

Înainte de a da exemple de instrucțiuni grafice, câteva cuvinte despre culori. Culorile implicite pe care le utilizează R la un moment dat pentru reprezentări grafice sunt stocate într-o paletă, un vector de tip caracter cu numele culorilor. Le putem vedea cu comanda *palette()*.

```
> palette()
[1] "black" "red" "green3" "blue" "cyan" "magenta"
[7] "yellow" "gray"
```

Pentru culorile care nu au nume se utilizează descrierea lor hexazecimală. Putem crea palete de culori prin *palette(rainbow(n))*, *palette(heat.colors(n))*, *palette(terrain.colors(n))*, *palette(topo.colors(n))*, *palette(cm.colors(n))*, unde *n* este numărul culorilor pe care le vrem în paletă. O altă modalitate de a obține palete este prin interpolarea culorilor cu funcția *colorRampPalette* ca mai jos:

```
> culorile.mele=colorRampPalette(c("blue","yellow","red"))
> culorile.mele
function (n)
{
  x <- ramp(seq.int(0, 1, length.out = n))
  rgb(x[, 1], x[, 2], x[, 3], maxColorValue = 255)
}
<bytecode: 0x01aa3ae0>
<environment: 0x02c36ab8>
> paleta=palette(culorile.mele(20))
> paleta
[1] "blue"      "#1A1AE4"  "#3535C9"  "#5050AE"  "#6B6B93"
[6] "#868678"  "#A1A15D"  "#BBBB43"  "#D6D628"  "#F1F10D"
[11] "#FFD100"  "#FFD600"  "#FFBB00"  "#FFA100"  "#FF8600"
[16] "#FF6B00"  "#FF5000"  "#FF3500"  "#FF1A00"  "red"
```

Se vede că funcția *colorRampPalette* ne întoarce ca rezultat o funcție de calcul a culorilor. Această funcție aplicată unui număr natural *n* crează o paletă de lungime *n*. Putem vedea culorile din paletă tipărind de exemplu:

```
> plot(1:20, col=paleta, pch=1:20, main="Puncte si culori")
```

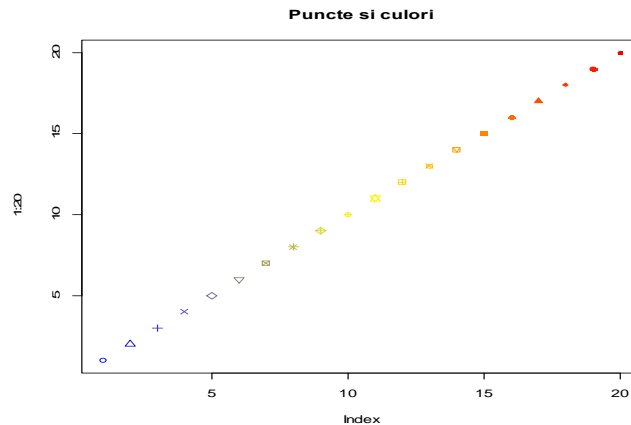


Fig 4.1

Un clic dreapta pe grafic face să apară un meniu contextual unde apar opțiunile de salvare a graficului în format metafile sau postscript sau de copiere în clipboard în format metafile sau bitmap (în windows). De aici este ușor să-l includem într-un document.

Pe lângă parametrii din vectorul *par* în comenzile grafice de nivel înalt apar și alți parametri opționali:

Parametrul	Explicații
add	Dacă este TRUE atunci noul grafic se suprapune peste cel existent. Funcționează uneori.
axes	Dacă este FALSE nu se trasează axele
log="x", log="y", log="xy"	Se specifică ce axe sunt gradate logaritmice. Acest lucru se poate face și prin parametrii grafici, <i>xlog</i> , <i>ylog</i> .
type	Poate lua valorile "p" pentru a tipări puncte, "l" pentru a tipări linii, "b" pentru a tipări puncte și linii între ele, "o" pentru a tipări puncte și linii peste ele, "s" sau "S" pentru a tipări grafice sub formă pas (step) în sus sau jos, "n" pentru a se seta sistemul de coordonate, fără a trasa graficul
xlab, ylab	Sunt șiruri de caractere pentru etichetarea axelor
main, sub	Sunt șiruri de caractere pentru titlu și subtitlu
xlim, ylim, zlim	Sunt vectori numerici cu câte două valori pentru a seta limitele pe cele trei axe

În cele ce urmează exemplificăm utilizarea în parte a rutinelor grafice, cu comentarii.

Exemplul 4.1. Facem graficele a trei funcții, conform scriptului următor:

```
x=seq(-pi,pi,length=20)
y1=sin(x); y2=cos(x); y3=2*x/(1+x^2)
plot(x,y1,type="b", pch="+", xlab="", ylab="", lwd=3)
lines(x,y2, type="o", pch=10, col="red", lty="dotted",
lwd=3)
lines(x,y3,type="l", lty="dashed", col="blue", lwd=3)
title(main="Un exemplu de grafic", sub="Graficele a trei
functii",xlab="axa X",ylab="axa Y", cex.main=2,
cex.lab=1.5, cex.sub=2)
grid(5,5,col="cyan")
legend("topleft", legend=c("sin","cos","alta")
,col=c("black","red","blue"),
lty=c("solid","dotted","dashed")
,bg="yellow1")
```

Rezultatul este în Fig 4.2.

Observații:

a) Este bine când se reprezintă grafic mai multe funcții, la apelul rutinei plot să se specifice $xlim=range(x1,x2,x3,...)$, $ylim=range(y1,y2,y3,...)$ unde $x1$, $x2$, $x3$, .. sunt vectorii cu abscisele coordonatelor punctelor de pe grafic iar $y1$, $y2$, $y3$, .. conțin ordonatele punctelor. Acest lucru ne asigură că toate graficele vor fi complet reprezentate. Dacă la reprezentările din figura de mai sus mai adăugăm un grafic ce are puncte în afara zonei delimitate (în cazul nostru $[-3,3] \times [-1,1]$) acestea nu vor apare.

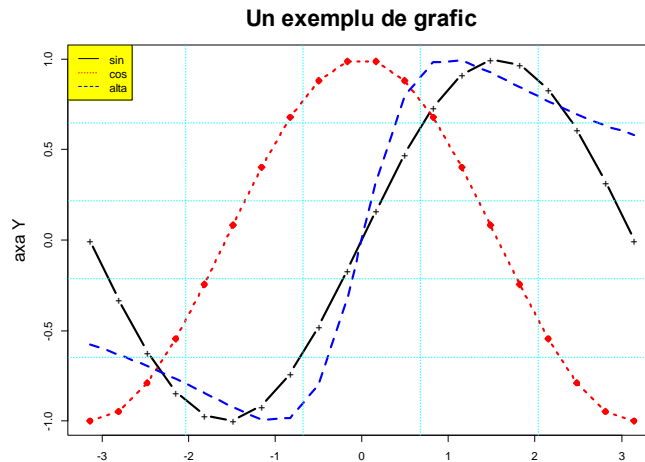
b) Primul grafic se execută cu *plot* iar celelalte cu *lines*, *points* sau alte comenzi de nivel inferior pentru a se adăuga elementele la ceea ce este deja reprezentat.

c) Titlurile se pot adăuga într-o instrucțiune *plot*, sau prin comanda *title*. Atenție la faptul că adăugarea ulterioară a unui titlu nu schimbă altul existent ci se suprapune celui existent.

d) Poziționarea legendei se poate face prin cuvintele "topleft", "topright", "top", "bottomleft", "bottomright", "bottom", "left", "right", sau prin coordonatele x , y

ale colțului stânga sus, coordonate ce pot fi alese interactiv prin *locator()*. Instrucțiunea *legend* din exemplul precedent ar fi putut fi

```
legend(locator(), legend=c("sin", "cos", "alta"),
col=c("black", "red", "blue"),
lty=c("solid", "dotted", "dashed"),
bg="yellow1")
```



Graficele a trei functii
Fig 4.2

Exemplul 4.2. Graficul unei suprafețe văzută din diverse direcții. Graficele sunt grupate într-o singură figură divizată în 2x2 zone prin comanda *par(mfrow(2,2))*.

```
f<-function(x,y){z=x-sin(2*x)^2+2*y^2}
x=seq(-1,1,length=20); y=x;
z=outer(x,y,f);

f<-function(x,y){z=x-sin(2*x)^2+2*y^2}
x=seq(-1,1,length=20); y=x;
z=outer(x,y,f);

par(mfrow=c(2,2))

persp(x,y,z,col="lightblue", shade=0.5, theta=45, phi=30,
r=50,
xlim=1.1*range(x),ylim=1.1*range(y),zlim=1.1*range(z),
main="Primul grafic", sub="Graficul functiei z=x-
sin(2*x)^2+2*y^2",
xlab="axa X",ylab="axa Y", zlab="Z",
ticktype="detailed", cex.main=2, cex.sub=2)

paleta=palette(rainbow(32))
clase=cut(z,32)
par(mfg=c(1,2))
persp(x,y,z,col=paleta[clase], shade=0.5, theta=30,
phi=15, r=50,
xlim=1.1*range(x),ylim=1.1*range(y),zlim=1.1*range(z),
```

```

main="Al doilea grafic", sub="Graficul functiei  $z=x-\sin(2*x)^2+2*y^2$ ",
xlab="axa X",ylab="axa Y", zlab="Z",
ticktype="detailed",cex.main=2, cex.sub=2)

paleta1=palette(cm.colors(32))
clase=cut(z,32)
persp(x,y,z,col=paleta1[clase], shade=0.5, theta=20,
phi=45, r=50,
xlim=1.1*range(x),ylim=1.1*range(y),zlim=1.1*range(z),
main="Al treilea grafic", sub="Graficul functiei  $z=x-\sin(2*x)^2+2*y^2$ ",
xlab="axa X",ylab="axa Y", zlab="Z",
ticktype="detailed",cex.main=2, cex.sub=2)

paleta2=palette(colorRampPalette(c("maroon", "gold"))(32))
clase=cut(z,32)
persp(x,y,z,col=paleta2[clase], shade=1, theta=30, phi=45,
r=50,
xlim=1.1*range(x),ylim=1.1*range(y),zlim=1.1*range(z),
main="Al patrulea grafic", sub="Graficul functiei  $z=x-\sin(2*x)^2+2*y^2$ ",
xlab="axa X",ylab="axa Y", zlab="Z",
ticktype="detailed",cex.main=2, cex.sub=2)

par(mfrow=c(1,1))

```

Rezultatul este în Fig 4.3. Se constată că nu se respectă paletele de culori utilizate la fiecare figură. Când se utilizează o singură culoare pentru suprafață ea este imediat luată imediat în considerație pe când dacă se variază culoarea de la o fațetă la alta atunci paletele utilizate se permută cumva între grafice.

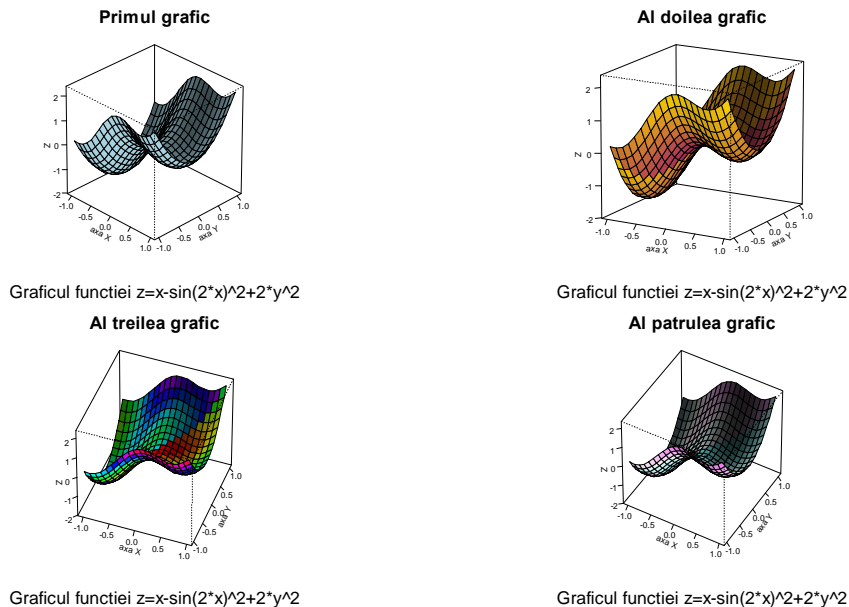


Fig 4.3

Exemplul 4.3. Pentru adnotări matematice pe grafice se pot obține informații din help prin comanda *help(plotmath)*. Expresiile matematice se alcătuiesc cu comanda *expression(...)* în care se introduce sub forma unei expresii apropiate de cea latex formula matematică dorită. Pentru a lipi între ele mai multe șiruri se utilizează comanda *paste(sir1, sir2,...)*. În tabelul de mai jos avem câteva astfel de expresii.

Expresia în R	Forma pe grafic după executia comenzii <i>text(locator(), ex)</i>
<code>ex=expression(paste(x*y, " ", "x%*%y", " ", x/y, " ", "x%//%y", " ", "x%+-%y"))</code>	$xy, x \times y, x/y, x \div y, x \pm y$
<code>ex=expression(paste(x^2, " ", x[2], " ", sqrt(x), " ", sqrt(x,3)))</code>	$x^2, x_2, \sqrt{x}, \sqrt[3]{x}$
<code>ex=expression(paste(x == y, " ", x!=y, " ", x<=y, " ", x %~% y, " ", x %==% y))</code>	$x = y, x \neq y, x \leq y, x \approx y, x \equiv y$
<code>ex=expression(paste(cdots, " ", x %subset% y, " ", ldots, " ", x %notsubset% y, " ", x %in% y, " ", x %notin% y))</code>	$\cdots, x \subset y, \dots, x \not\subset y, x \in y, x \notin y$
<code>ex=expression(paste(hat(x), " ", tilde(x), " ", dot(x), " ", ring(x), " ", bar(xyzu), " ", widehat(xyzu)))</code>	$\hat{x}, \tilde{x}, \dot{x}, \overset{\circ}{x}, \overline{xyzu}, \widehat{xyzu}$
<code>ex=expression(paste(x %>% y, " ", x %up% y, " ", x %=>% y, " ", alpha, " ", Alpha, " ", infinity, " ", nabla))</code>	$x \rightarrow y, x \Uparrow y, x \Rightarrow y, \alpha, A, \infty, \nabla$
<code>ex=expression(paste(underline(x), " ", frac(x, y), " ", sum(x[i], i==1, n), " ", integral(f(x)*dx, a, b), " ", union(A[i], i==1, n), " ", lim(f(x), x %>% 0)))</code>	$\underline{x}, \frac{x}{y}, \sum_{i=1}^n x_i, \int_a^b f(x)dx, \bigcup_{i=1}^n A_i, \lim_{x \rightarrow 0} f(x)$
<code>ex=expression(paste(y=sqrt(frac(x^2-1,x^2+1))%*%sin(x)))</code>	$\sqrt{\frac{x^2-1}{x^2+1}} \times \sin(x)$
<code>ex=expression(paste(bgroup((" ", frac(x+y,x-y)+frac(x^2+y^2,x^2-y^2),"))^2))</code>	$\left(\frac{x+y}{x-y} + \frac{x^2+y^2}{x^2-y^2} \right)^2$

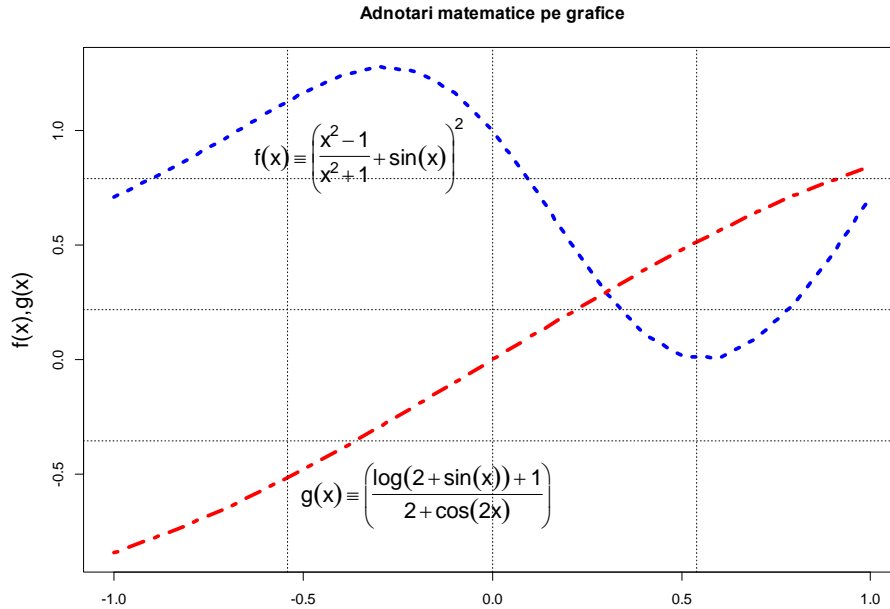
Utilizarea comenzii *text(locator(), ex)* pentru a pune expresia pe grafic se justifică prin aceea că *locator()* așteaptă un clic stânga pe o poziție din grafic pentru a se stabili unde va fi pus textul. După aceea se face un clic dreapta și din meniul derulant se alege opțiunea stop. Mai jos avem două grafice cu adnotări.

```
f<-function(x) ((x^2-1)/(x^2+1)+sin(x))^2
g<-function(x) (ln(2+sin(x))+1)/(2+cos(2*x))
x=seq(-1,1,length=21); y1=f(x); y2=sin(x);
plot(x,y1,ylim=range(y1,y2),type="l", lty=3, col="blue",
lwd=4, xlab="x",
ylab="f(x), g(x)", cex.lab=1.5)
lines(x,y2,type="l", lty=4, col="red",lwd=4)
grid(4,4,col="black")
title(main="Adnotari matematice pe grafice",
sub="Graficul a doua functii", cex.sub=1.5)
```

```

ex1=expression(paste(f(x)===% bgroup("(",frac(x^2-
1,x^2+1)+sin(x),")")^2))
text(locator(),ex1,cex=1.5)
ex2=expression(paste(g(x)===%
bgroup("(",frac(log(2+sin(x))+1,2+cos(2*x)),")"))))
text(locator(),ex2, cex=1.5)

```



Graficul a două funcții
Fig 4.4

Exemplul 4.4. Histograma valorilor generate aleator după o distribuție normală este comparată cu graficul distribuției.

```

media=2; devs=4;
y=rnorm(nr, mean=2, sd=4)
hist(y, probability=TRUE, col="lightgray",
main="Histograma valorilor y=rnorm(200, mean=2, sd=4)",
xlab="Valorile lui y",
ylab="Densitatea valorilor lui y", cex.axis=1.5,
cex.lab=1.5 )
x1=seq(-6,10,length=100)
y1=dnorm(x1,mean=2, sd=4)
lines(x1,y1, lwd=3)
ex=expression(paste(rho, "(", x, ")=",frac(1,
4*sqrt(2*pi)), " ",
plain(e)^{frac(-(x-2)^2, 2*%4^2)}))
text(locator(),ex, cex=1.5)

```

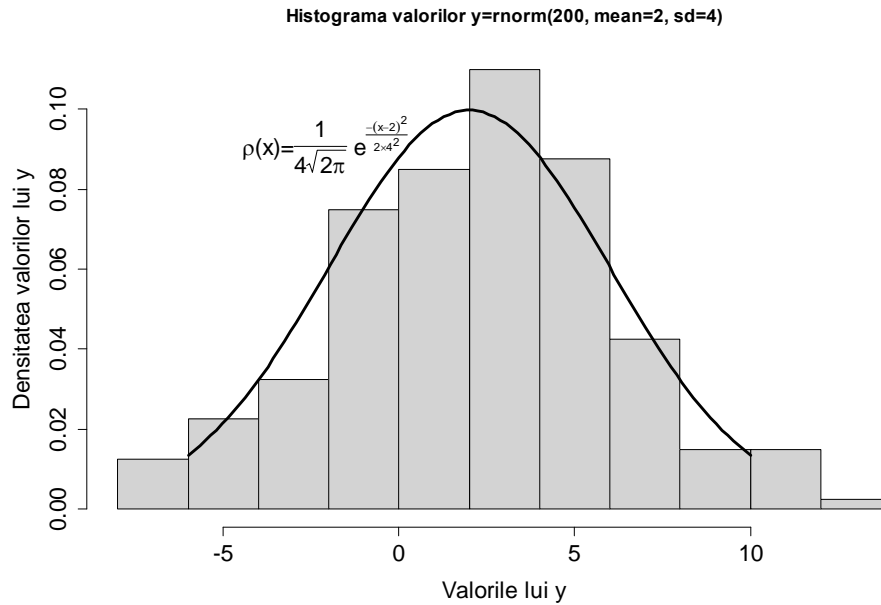



Fig 4.5

Exemplul 4.5. Utilizăm *qqnorm* pentru a detecta dacă o selecție este normală (graficul ar trebui să fie o dreaptă) și *qqplot* pentru a detecta dacă două selecții empirice au aceeași distribuție (graficul ar trebui să fie prima bisectoare).

```
x1=runif(1000);x2=rnorm(1000, mean=1, sd=1)
par(mfrow=c(1,2));
qqnorm(x2, main="Q-Q norm pentru o selectie normala",
xlab="Cuantilele teoretice",
ylab="x2=rnorm(100, mean=1, sd=1)", pch="+");
qqplot(x1,x2, main="Q-Q plot pentru doua selectii",
pch="+",
xlab="x1=runif(100)", ylab="x2=rnorm(100, mean=1, sd=1)");
par(mfrow=c(1,1))
```

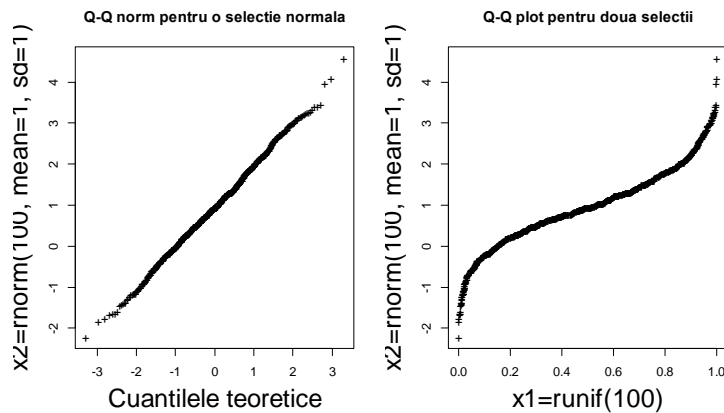


Fig 4.6

Capitolul 5

Programarea în R

În principiu, un limbaj de programare ne pune la dispoziție un anumit număr de tipuri de date și un anumit număr de operații predefinite asupra lor, modalități de introducere a datelor într-un program precum și de extragere a rezultatelor (instrucțiuni read/write), instrucțiuni de control al fluxului de operații executate (în principal instrucțiuni de decizie și repetiție), instrucțiuni de creare de subprograme (funcții). În funcție de limbaj există și alte aspecte importante care pot interesa: posibilitatea de a utiliza rutine scrise în alte limbaje, posibilitatea de utilizare recursivă a funcțiilor, programarea pe obiecte, etc. În cele ce urmează descriem succint elementele de bază ale programării în limbajul R.

Codul scris în limbajul R este cod interpretat, deci instrucțiunile sunt compilate la un stadiu intermediar (bytecode) de unde sunt transformate în cod dependent de mașină în momentul execuției. Acest proces de interpretare face ca timpul de execuție să fie mai lung decât pentru programele compilate de la început în cod-mașină executabil.

Scrierea unui program R se poate face în orice editor text. Se poate utiliza editorul propriu (apelabil cu opțiunea “New script” din meniul File) care are avantajul că din el se poate rula o porțiune din program selectată în prealabil printr-un clic dreapta și alegerea opțiunii “Run line or selection” din meniul de context. Dacă programul a fost salvat pe disc în fișierul prog.r atunci rularea lui se face cu comanda source(“prog.r”) dacă este în directorul curent, altfel trebuie indicată și calea.

Toate instrucțiunile care se scriu la consolă într-o sesiune de lucru pot fi salvate la final în fișierul text .Rhistory (“Save history” din meniul File) de unde se pot rula din nou cu comanda “Load history” (din meniul File). În acest fel se poate repeta prelucrarea unor date într-o sesiune nouă sau se poate scrie un program după ce se execută pas cu pas instrucțiunile lui ca și comenzi la consolă.

5.1. Tipurile de date în R. Citirea și scrierea lor. Operații predefinite

Capitolul 2 este dedicat tipurilor de date, operațiilor asupra lor și modalităților de a le crea, citi sau extrage din mediul R. În capitolul 3 se da o listă mai extinsă a operațiilor predefinite asupra datelor recunoscute de R. Aceste operații sunt prezentate sub forma unor funcții. Aici menționăm doar că datele simple cu care lucrează R sunt alcătuite din numere, valori booleene, caractere (șiruri de caractere), date binare (raw). Tipurile fundamentale de date compuse din date simple sunt vectorii, matricele, blocurile multidimensionale (array), data.frame, listele. Se utilizează de asemenea NA pentru date indisponibile și NaN pentru valori nenumarice. Pot fi create tipuri noi de date.

Un lucru important care trebuie avut în vedere este posibilitatea de grupare a instrucțiunilor prin acolade. Astfel {instr1; instr2; instr3} este un grup. Aceste grupuri apar în general în instrucțiuni de decizie sau repetiție. Dacă sunt pe același rând instrucțiunile se separă prin “;” iar dacă sunt pe rânduri diferite nu e nevoie de vreun semn de separație.

Pentru comentarea unor secțiuni de program se utilizează caracterul # la începutul liniei de comentariu.

După efectuare unui calcul dacă vrem să păstrăm rezultatele ele trebuie atribuite unei variabile. Atribuirea se face prin (variabila care primește valoarea atribuită este A):

```
A<-expresie
A<<-expresie
A=expresie
expresie->A
expresie->>A
```

Atribuirile <- sau -> sunt locale iar <<- și ->> sunt globale (vezi observația f) de la secțiunea “definirea funcțiilor proprii”). Atribuirea prin = este echivalentă cu atribuirea -> (există unele situații când nu sunt echivalente dar nu le utilizăm)

5.2. Instrucțiunea de decizie

Instrucțiunea de decizie este IF. Forma ei este la fel ca în alte limbaje:

```
if (conditie) instr_1 else instr_2
```

- i. `conditie` este o expresie care dă o valoare logică (TRUE, FALSE)
- ii. `instr_1` este o instrucțiune sau un grup de instrucțiuni și la fel `instr_2`
- iii. `else` nu poate să apară singur pe o linie. Este de preferat să fie pus după instrucțiunea (sau grupul) `instr_1`.

Exemplul 5.1.

```
a=5; b=3
if ( (a==5) & (b!=a) )
{ s<-a+b^2
  s<-s^2 } else
{ s=a-b^2; s<-s^2 }
s
[1] 196
```

Există și o funcție `ifelse` de forma:

```
ifelse (conditie, da, nu)
```

Dacă este adevărată afirmația *conditie* atunci rezultatul lui *ifelse* este *da* altfel rezultatul este *nu*. În *conditie* putem avea un vector care se evaluează pe fiecare componentă la TRUE sau FALSE. Rezultatul funcției este un vector de aceeași lungime ca vectorul *conditie*, cu valori din vectorii *da* sau *nu*. Dacă vectorii *da* sau *nu* nu au aceeași lungime ca și vectorul *conditie* atunci sunt trunchiați sau extinși prin repetare.

Exemplul 5.2.

```
> a=c(1,-2,3,0,-4,-5,0,2)
> da=c(1,2,3)
> nu=c(-1,-2,-3,-4,-5,-6,-7,-8,-9,-10)
> ifelse(a>0,da,nu)
[1] 1 -2 3 -4 -5 -6 -7 2
```

3. Instrucțiuni de repetiție

Instrucțiunile de repetiție (ciclare) sunt în R: *for*, *while*, *repeat*. Aceste instrucțiuni au formele:

for (*var1 in var2*) *instr*

unde *var2* este o variabilă ce se poate converti la un șir de valori (ca de exemplu vector sau listă), iar *instr* este o instrucțiune sau un grup de instrucțiuni.

while (*cond*) *instr*

unde *cond* este o expresie ce se evaluează la TRUE sau FALSE, iar *instr* este o instrucțiune sau un grup de instrucțiuni. Se va avea desigur grijă ca variabila *cond* să se modifice în grupul de instrucțiuni *instr* pentru a deveni la un moment FALSE și a se opri repetiția.

repeat *instr*

este o variantă de repetiție în care *instr* este un bloc de instrucțiuni. Pentru a se asigura că se termină la un moment aceste repetiții trebuie întreținută în blocul *instr* o condiție care se evaluează la TRUE sau FALSE. Ieșirea din bloc se face la execuția unei instrucțiuni de forma: *if*(condiție) *break*, dacă condiție=TRUE.

În blocurile de repetiție se poate folosi comanda ***break*** ce determină ieșirea imediată din ciclu sau comanda ***next*** care determină trecerea la următoarea repetiție din ciclu. În cazul ciclurilor incluse unul în altul *break* și *next* se utilizează doar pentru ciclul cel mai interior.

Exemplul 5.3. Se dă o matrice *a* de dimensiune 3x4 și se cere să se calculeze vectorul *s* de dimensiune egală 3 pentru care *s*[*i*] este suma elementelor din linia *i*. Rezolvare prin repetiție cu *for*:

```
a=matrix(c(1:12),nrow=3)
s=rep(0,3)
for (i in 1:3){
  for (j in 1:4) s[i]<-s[i]+a[i,j]
}

> s
[1] 22 26 30
```

Rezolvare prin repetiție cu *while*

```
a=matrix(c(1:12),nrow=3)
s=rep(0,3)
i=1
while(i<=3){
  j=1
  while (j<=4) {s[i]=s[i]+a[i,j];j=j+1}
  i=i+1
}

> s
[1] 22 26 30
```

Rezolvare prin repetiție cu *repeat*

```

a=matrix(c(1:12),nrow=3)
s=rep(0,3)
i=1
repeat{
  j=1
  repeat {s[i]=s[i]+a[i,j];j=j+1; if(j==5) break}
  i=i+1
  if(i==4) break
}

> s
[1] 22 26 30

```

În toate aceste programe ciclul interior se poate înlocui cu funcția *sum*. Prin aceasta se execută mai repede programul pentru că repetiția din funcția *sum* este compilată direct în cod mașină. Programul ar arăta astfel (în varianta cu *for*).

```

a=matrix(c(1:12),nrow=3)
for (i in 1:3) s[i]=sum(a[i,])

> s
[1] 22 26 30

```

Multe repetiții pot fi realizate cu funcțiile de mapare *apply*, *lapply*, *sapply*, *tapply*, *mapply*, descrise în capitolul 3. Programul precedent poate fi scris de exemplu cu funcția *apply* astfel:

```

a=matrix(c(1:12),nrow=3)
s=apply(a,MARGIN=1,sum)

> s
[1] 22 26 30

```

Exemplul 5.4. Date matricele a , b , c de dimensiune 3×4 să se determine matricea d de dimensiune 3×4 pentru care $d[i, j]$ este media a vectorului format de $a[i, j]$, $b[i, j]$, $c[i, j]$.

Repetiția cu *for* ne dă două cicluri:

```

a=matrix(1:12,nrow=3);
b=matrix(seq(1,2,length=12),nrow=3)
c=matrix(11:22, nrow=3)

d=matrix(rep(0,12),nrow=3)
for (i in 1:3)
  for (j in 1:4)
    d[i,j]=mean(a[i,j],b[i,j],c[i,j])

> d
      [,1] [,2] [,3] [,4]
[1,]    1    4    7   10
[2,]    2    5    8   11
[3,]    3    6    9   12

```

Repetiția prin funcția *mapply* se realizează astfel:

```
a=matrix(1:12,nrow=3);
b=matrix(seq(1,2,length=12),nrow=3)
c=matrix(11:22, nrow=3)

d=mapply(mean,a,b,c)
d=matrix(d,nrow=3)

> d
      [,1] [,2] [,3] [,4]
[1,]     1     4     7    10
[2,]     2     5     8    11
[3,]     3     6     9    12
```

5.4. Definirea funcțiilor proprii

În limbajul R definirea funcțiilor se face sub forma

nume_functie<-function(argumente) {instructiuni de calcul}

iar apelul funcției se face simplu *nume_functie(argumente)*. Rezultatul întors de funcție trebuie să apară ultimul, fără atribuire. O altă posibilitate este instrucțiunea *return(val)*. În acest caz, execuția instrucțiunilor din funcție se încheie imediat și *val* este returnată ca valoare a funcției.

Exemplul 5.5. Să definim funcția $f(x, y, z) = x^2 + xyz$ și apoi să o aplicăm argumentelor $x=3, y=5, z=9$.

```
f<-function(x,y,z){r<-x^2+x*y*z;r}
f(3,5,9)

[1] 144
```

În continuare am strâns câteva observații referitoare la funcții

a) Deoarece în definiția funcției am folosit numele x, y, z , dacă la apelul funcției vom utiliza argumentele cu numele lor, atunci ordinea în care apar poate fi schimbată.

```
> f<-function(x,y,z){r<-x^2+x*y*z;r}
> f(3,5,9)
[1] 144
> f(y=5,z=9,x=3)
[1] 144
> f(3, z=9,y=5)
[1] 144
> f(3,y=5,9)
[1] 144
> f(3,9,y=5)
[1] 144
```

b) Dacă unele argumente sunt prezente cu numele lor și altele fără nume, atunci cele fără nume trebuie să apară în ordinea în care apar în definiția funcției.

c) Este posibil ca în definiția unei funcții unele argumente să aibă valori implicite prestabilite. În acest caz la utilizarea funcției nu este neapărat necesar ca ele să fie specificate decât dacă se dorește utilizarea lor cu alte valori decât cele implicite. Mai jos avem un exemplu:

```
> f<-function(x,y=5,z){r<-x^2+x*y*z;r}
> f(3,5,9)
[1] 144
> f(x=3, z=9)
[1] 144
```

Un apel $f(3,9)$ pentru funcția precedentă generează un mesaj de eroare pentru că se face identificarea $x<-3, y<-9$ și nu este definită valoarea pentru z .

```
> f(3,9)
Error in x * y * z : 'z' is missing
```

d) Putem utiliza în definiția unei funcții și alte funcții definite anterior. De asemenea putem utiliza în definiția unei funcții argumentul “...” (trei puncte) care reprezintă argumente opționale. De obicei aceste argumente sunt transferate funcției interioare.

```
> g<-function(a=4,b=6) a+b/2
>
> f<-function(x,y=5,z,...){r1<-x^2+x*y*z; r2<-g(...)
+ r<-r1+r2; r}
>
> f(1,2,3)
[1] 14
> f(1,2,3,4); # a=4 ; b=6 implicit
[1] 14
> f(z=3,x=1,y=2,4,6); #a=4, b=6
[1] 14
> f(z=3,x=1,y=2,4,8); #a=4, b=8
[1] 15
> f(z=3,x=1,y=2,b=8,a=4); #a=4, b=8
[1] 15
> f(z=3,x=1,y=2,b=8); #a=4 (implicit), b=8
[1] 15
```

e) În anumite cazuri avem nevoie de o funcție doar într-un loc. În acest caz funcția se poate defini anonim (fără nume) și în consecință poate fi utilizată doar în acel loc. De exemplu să calculăm pentru trei matrice a, b, c de aceleași dimensiuni expresia $a^2+a*b*c$ (nu uităm că operațiile sunt vectorizate, se aplică elementelor în poziții similare în matricele a, b, c). Scriptul următor realizează acest lucru în două feluri, în unul din feluri se utilizează o funcție anonimă.

```
a=matrix(1:12, nrow=3)
b=matrix(seq(from=1,by=2, length=12), nrow=3)
```

```

c=matrix(rep(2,12), nrow=3)
rez1=a^2+a*b*c
rez2=matrix(mapply(function(x,y,z)x^2+x*y*z,a,b,c),nrow=3)
rez1
rez2

```

Rezultatele sunt aceleași pentru cele două metode:

```

> rez1
      [,1] [,2] [,3] [,4]
[1,]     3    72   231   480
[2,]    16   115   304   583
[3,]    39   168   387   696
> rez2
      [,1] [,2] [,3] [,4]
[1,]     3    72   231   480
[2,]    16   115   304   583
[3,]    39   168   387   696

```

f) Variabilele definite în corpul unei funcții, să zicem f , au caracter local. Ele sunt vizibile doar în corpul funcției f . De asemenea în corpul funcției f se pot utiliza date sau funcții definite în funcția în care este definită f , sau date ori funcții globale. Atunci când se caută valoarea unei variabile (sau definiția unei funcții) se caută printre parametrii locali (inclusiv argumentele funcției), apoi printre parametrii locali funcției superioare (în interiorul căreia este definită) s.a.m.d. Modificările parametrilor locali nu modifică parametrii globali cu același nume. Pentru a modifica în corpul unei funcții un parametru global, să zicem A , trebuie să utilizăm instrucțiunea de atribuire $A <- expresie$. La crearea unei funcții se crează un context de vizibilitate (environment) cu ceea ce este vizibil de către funcția respectivă.

Comentăm acum (și modificăm) un exemplu din manual: *An introduction to R*, pg. 47 (manualul vine cu distribuția standard de R).

Apelul funcției `open.account` produce o listă ce are ca elemente trei funcții (`deposit`, `withdraw`, `balance`). Crearea unui cont se face printr-o instrucțiune de tipul `client=open.account(sumadepusa)`. Variabila `data` este produsă de program. Variabila `dob` reprezintă rata anuală a dobânzii (compuse) și poate fi modificată de utilizatorul scriptului: este aceeași pentru orice client

<pre> open.account <- function(total=0,data,dob=0.10){ data=Sys.time() list(deposit = function(amount) { </pre>	<p>La crearea unui nou cont se apelează o singură dată funcția sub forma <code>client=open.account(...)</code></p> <p>Sistemul generează un context de vizibilitate pentru rezultatul funcției (în cazul nostru lista <code>client</code>). Funcțiile din lista <code>client</code> au în contextul lor de vizibilitate parametrii de intrare ai funcției generatoare <code>open.account</code>. În acest context vor fi văzute variabilele <code>total</code>, <code>dob</code>, <code>data</code>; <code>data</code> este data creării contului</p> <p>Prima funcție: <code>deposit</code> modifică</p>
--	---

<pre> if(amount <= 0) stop("Deposits must be positive!\n") dataa=Sys.time() durata=as.numeric(dataa) - as.numeric(data) durata=durata/86400 durata=durata/365 total <- total*(1+dob)^durata + amount cat(amount, "deposited. Your balance is", total, "\n\n") data<-dataa }, withdraw <- function(amount) { dataa=Sys.time() durata<-as.numeric(dataa) - as.numeric(data) durata<-durata/86400 durata<-durata/365 total1<-total*(1+dob)^durata if(amount > total1) stop("You don't have that much money!\n") total <- total1 - amount data<-dataa cat(amount, "withdrawn. Your balance is", total, "\n\n") }, balance <- function() { dataa=Sys.time() durata=as.numeric(dataa) - as.numeric(data) durata=durata/86400 durata=durata/365 total1=total*(1+dob)^durata cat("Your balance is", total1, "last change ",as.character(data),"\n\n") }) } </pre>	<p>valoarea contului ținând seama de dobîndă și valoarea depusă dataa este data alimentării contului Se convertesc datele în secunde, apoi în zile, apoi în ani pentru a se calcula durata de la ultima operație în cont Se actualizează suma după formula $total = total(1+dob)^{durata} + amount$ Se afișează informațiile pentru client Actualizează data</p> <p>A doua funcție: withdraw modifică valoarea contului ținând seama de dobîndă și valoarea retrasă</p> <p>Se actualizează totalul după formula $Total1 = total(1+dob)^{durata}$ Dacă total1 este prea mic, se dă mesaj de eroare Dacă nu, se modifică valoarea contului Actualizează data Mesaj de confirmare</p> <p>A treia funcție: balance nu modifică nimic în cont, nici data, doar informează clientul asupra sumei disponibile și asupra datei ultimei operații în cont.</p>
---	--

După ce se rulează scriptul ce conține funcția open.account se poate lucra în flul următor cu conturile:

```

> gh<-open.account(10000000)
> gh$deposit(1000000)
1e+06 deposited. Your balance is 1.1e+07

> gh$withdraw(100000)
1e+05 withdrawn. Your balance is 10900001

> gh$balance()
Your balance is 10900002 last change 2011-12-19 01:23:47

```

Se vede cum cu trecerea timpului se adaugă dobânda.

g) Printre parametrii unei funcții pot fi alte funcții. De asemenea în corpul unei funcții pot fi apeluri la ea însăși ceea ce duce la aplicarea recursivă a funcției. În acest caz e nevoie de condiții care să asigure că numărul de aplicări recursive e finit. Ca exemplu avem următorul script simplu ce calculează recursiv termenii șirului Fibonacci.

```
fib<-function(n){
  if(n==1) return(c(1))
  if(n==2) return(c(1,1))
  s1=fib(n-1)
  urm=s1[n-1]+s1[n-2]
  return (c(s1,urm))
}
> fib(10)
[1] 1 1 2 3 5 8 13 21 34 55
```

5.5. Apelul din R al funcțiilor scrise în C

Apar situații când nu avem în R rutine pentru anumite operații. Soluția este scrierea de rutine (funcții) proprii. Utilizarea într-o măsură mare a instrucțiunilor de decizie sau repetiție în R face ca timpul de rulare să crească mult datorită faptului că R este un limbaj interpretat. O ieșire o reprezintă uneori utilizarea de rutine scrise într-un limbaj în care se face compilarea la nivel de cod mașină, cum ar fi C, Fortran, Pascal, etc. În R există mecanisme de apel a funcțiilor scrise în C ori Fortran precum și mecanisme de apel într-un program extern a unor rutine din R. Ne vom ocupa de metoda de apel în R a unor rutine scrise în C.

Atunci când se apelează într-un program funcții scrise în alt program trebuie avute în vedere corespondența tipurilor de date, modul de transfer al datelor spre rutina apelată, modul de preluare a rezultatelor produse de rutina apelată. În documentația standard pentru R găsim următoarea corespondență între tipurile de date din R, C, Fortran :

R storage mode	C type	FORTRAN type
logical	int *	INTEGER
integer	int *	INTEGER
double	double *	DOUBLE PRECISION
complex	Rcomplex *	DOUBLE COMPLEX
character	char **	CHARACTER*255
raw	unsigned char *	none

Pentru cei care doresc să scrie pachete de programe pentru R utilizând limbajul C, în subdirectorul « include » din directorul unde este instalat R se găsesc fișiere de definiție pentru diverse tipuri de date și diverse funcții din librării larg utilizate. Pentru noi corespondența de tipuri din tabelul de mai sus este suficientă.

Câteva observații sunt de avut în vedere:

a) Datele între R și rutine externe se schimbă doar prin adrese.

b) Rutinele externe nu trebuie să întoarcă vreo valoare (în C ele trebuie să fie funcții de tip void iar în Fortran trebuie să fie subprograme de tip subroutine)

c) Schimbul de date (parametri de intrare și rezultate) se face doar prin argumentele funcțiilor (subrutinelor în Fortran)

d) Rutinele externe (în C ori Fortran) se pun în biblioteci dinamice (.dll în windows, .so în linux) care se încarcă în R cu comanda *dyn.load(ume_librarie,)* și se deconectează de la R prin comanda *dyn.unload(ume_librarie)*.

e) În principiu se pot utiliza în R și rutine scrise în Pascal dacă acestea sunt declarate de tip C. Declararea în Pascal a unor rutine ca fiind de tip C este importantă pentru a se face corect corespondența între parametri (în Pascal transmiterea parametrilor către o subrutină se face în ordine inversă față de C). De asemenea trebuie avut grijă la corespondența tipurilor de date (de exemplu double din R corespunde cu ^double în Pascal)

f) Apelul funcțiilor C se poate face cu comanda

.C(NAME, ..., NAOK = FALSE, DUP = TRUE, PACKAGE, ENCODING)

unde NAME este numele funcției, dacă NAOK este TRUE atunci valorile NA, NaN, Inf sunt transmise funcției C, altfel se generează o eroare dacă apar, DUP=TRUE determină o duplicare a datelor înainte de a fi transmise subrutinei, PACKAGE specifică în ce pachete să se caute subrutina NAME. Este obligatorie doar prezența parametrului NAME în apelul *.C(NAME,...)*

g) Pentru comoditate se poate introduce comanda *.C(NAME,...)* într-o funcție R care face apelul mult mai comod și se pot prelua din parametrii întorși doar pe cei care interesează (vezi exemplul).

În cele ce urmează ne mărginim la utilizarea unor rutine C care au ca parametri doar numere, reale sau întregi. Pentru aplicații matematice pare cea mai frecventă situație (ar mai fi eventual și utilizarea de numere complexe ca perechi de numere reale).

Exemplul 5.6. Sa se scrie o funcție C numită „rot” care plecând de la o matrice pătratică „a” să producă un vector care conține elementele matricei în ordinea din figura următoare, apoi să se apeleze această funcție din R.

$$\begin{pmatrix} 1 & 12 & 11 & 10 \\ 2 & 13 & 16 & 9 \\ 3 & 14 & 15 & 8 \\ 4 & 5 & 6 & 7 \end{pmatrix}$$

Ordinea de extragere a elementelor din matrice

Programul de extragere îl scriem în C. Am utilizat mediul de programare C++ gratuit Code:Blocks 8.02 (pentru windows) în care am ales un proiect nou de tip .dll numit pr1, care va avea ca efect crearea librăriei dinamice pr1.dll. În cadrul acestui proiect au fost create automat două fișiere main.h și main.dll de mai jos în care partea pe care am adăugat-o este evidențiată (se recunoaște rutina de extragere a elementelor din matricea a în vectorul b în ordinea specificată, precum și declarația că această rutină să aparțină librăriei numită aici pr1.dll)

Fișierul main.h

```

#ifndef __MAIN_H__
#define __MAIN_H__

#include <windows.h>

/* To use this exported function of dll, include this
   * header
   * in your project. */

#ifdef BUILD_DLL
#define DLL_EXPORT __declspec(dllexport)
#else
#define DLL_EXPORT __declspec(dllimport)
#endif

#ifdef __cplusplus
extern "C"
{
#endif

void DLL_EXPORT SomeFunction(const LPCSTR sometext);
void DLL_EXPORT rot(int *nlin, double *a, double *b);

#ifdef __cplusplus
}
#endif

#endif // __MAIN_H__

```

Fişierul main.cpp

```

#include "main.h"

// a sample exported function
void DLL_EXPORT SomeFunction(const LPCSTR sometext)
{
    MessageBoxA(0, sometext, "DLL Message", MB_OK |
MB_ICONINFORMATION);
}

void DLL_EXPORT rot(int *nlin, double *a, double *b){
    int i,j,k,cnt=0;
    int n=*nlin,mi=n/2;

    for (k=0;k<mi;k++){
        for(i=k;i<n-k-1;i++){b[cnt]=a[i+n*k];cnt++;}
        for(j=k;j<n-k-1;j++){b[cnt]=a[n-k-1+n*j];cnt++;}
        for(i=n-k-1;i>k;i--){b[cnt]=a[i+n*(n-k-1)];cnt++;}
        for(j=n-k-1;j>k;j--){b[cnt]=a[k+n*j];cnt++;}
    }
    if(n%2==1) b[cnt]=a[mi+n*mi];
}

BOOL WINAPI DllMain(HINSTANCE hinstDLL, DWORD fdwReason,
LPVOID lpvReserved)

```

```

{
    switch (fdwReason)
    {
        case DLL_PROCESS_ATTACH:
            // attach to process
            // return FALSE to fail DLL load
            break;

        case DLL_PROCESS_DETACH:
            // detach from process
            break;

        case DLL_THREAD_ATTACH:
            // attach to thread
            break;

        case DLL_THREAD_DETACH:
            // detach from thread
            break;
    }
    return TRUE; // succesful
}

```

După compilare și linkeditare (am optat pentru compilatorul GCC) se obține librăria `pr1.dll` pe care am mutat-o în directorul curent din R.

Programul R în care am utilizat funcția *rot* este următorul:

```

dyn.load("pr1.dll")
x=1:25
a=matrix(x,ncol=5)

b=rep(0,length(a))
rez<-
.C("rot",as.integer(ncol(a)),as.double(a),as.double(b))

dyn.unload("pr1.dll")

```

Rezultatele le vedem mai jos. Parametrii întorși sunt preluați de R sub forma unei liste în ordinea în care au fost declarați. Avem trei parametri în rutina *rot* deci și rezultatul întors *rez* va fi o listă cu trei intrări. Ultima intrare (a treia `rez[[3]]`) este vectorul *b* pe care îl căutăm.

```

> a
      [,1] [,2] [,3] [,4] [,5]
[1,]     1     6    11    16    21
[2,]     2     7    12    17    22
[3,]     3     8    13    18    23
[4,]     4     9    14    19    24
[5,]     5    10    15    20    25
> rez[[3]]
 [1]  1  2  3  4  5 10 15 20 25 24 23 22 21 16 11  6  7  8
 [9] 14 19 18 17 12 13

```

Pentru apelul funcției *rot* care se găsește în lbrăria *pr1.dll* putem scrie o funcție R care să facă acest apel mai ușor. În fișierul următor se face acest lucru. Din cei trei parametri îl luăm ca parametru de ieșire din funcția *rot* doar pe al treilea, b.

```
dyn.load("pr1.dll")
x=1:25
a=matrix(x,ncol=5)

rot<-function(a){
  if(!is.matrix(a))stop("intrarea in rot trebuie sa
fie matrice")
  if(!is.numeric(a)) stop("matricea trebuie sa fie
numERICA")
  if(!(ncol(a)==nrow(a))) stop("matricea trebuie sa
fie patratica")
  b=rep(0,length(a));
  rez<-
C("rot",as.integer(ncol(a)),as.double(a),as.double(b)
)
  return(rez[[3]])
}
aa<-rot(a)
dyn.unload("pr1.dll")
```

Rezultatul se vede mai jos:

```
> a
      [,1] [,2] [,3] [,4] [,5]
[1,]     1     6    11    16    21
[2,]     2     7    12    17    22
[3,]     3     8    13    18    23
[4,]     4     9    14    19    24
[5,]     5    10    15    20    25
> aa
[1]  1  2  3  4  5 10 15 20 25 24 23 22 21 16 11  6  7  8
9 14 19 18 17 12 13
>
```

Apelul rutinelor scrise în Fortran este asemanator. Ele sunt puse în librii dinamice de unde se încarcă cu comanda `dyn.load(ume_librarie)`. Apelul unei funcții din librărie se face prin instrucțiunea:

```
.Fortran(NAME, ..., NAOK = FALSE, DUP = TRUE, PACKAGE,
          ENCODING)
```

unde parametrii au aceeași semnificație ca la apelul funcțiilor C.

Capitolul 6

Aplicații diverse programate în R

Poate cel mai ușor mod de a învăța un limbaj este de a urmări câteva programe simple scrise în el.

6.1. Rambursări de credite.

Un client face un credit în valoare de S_0 UM la o bancă. În urma negocierilor, obține o rată a dobânzii egală cu i . Dorește să își facă un plan de rambursare în care să plătească ratele la momentele de timp t_1, \dots, t_n ; ar mai dori să știe și ce suma mai are de plătit după fiecare rată.

Să presupunem că ratele plătite la momentele t_i au valoarea r_i . Partea efectivă din datorie plătită atunci va fi egală cu $a_i = r_i - D_i$ unde D_i este dobânda plătită la suma restantă pe intervalul de timp (t_{i-1}, t_i) . Fie S_i suma restantă la momentul t_i , după achitarea ratei r_i .

Axioma este că dacă dobânda unitară pe intervalul (t_{i-1}, t_i) este d_i , atunci dobânda plătită la momentul t_i va fi egală cu $D_i = S_{i-1}d_i$. Să notăm $q_i = 1 + d_i$.

Cum prima dobândă se aplică la suma totală S_0 obținem următoarea recurență pentru sumele restante

$$\begin{aligned} a_1 &= S_0 - S_1 = r_1 - d_1 S_0 \Leftrightarrow S_1 = (1 + d_1)S_0 - r_1 \Leftrightarrow S_1 = q_1 S_0 - r_1 \\ a_2 &= S_1 - S_2 = r_2 - d_2 S_1 \Leftrightarrow S_2 = (1 + d_2)S_1 - r_2 \Leftrightarrow S_2 = q_2 S_1 - r_2 \\ &\dots\dots\dots \\ S_k &= (1 + d_k)S_{k-1} - r_k \Leftrightarrow S_k = q_k S_{k-1} - r_k \end{aligned} \quad (6.1.1)$$

Putem scrie relația (6.1.1) și sub forma $S_k + r_k = q_k S_{k-1}$. Să notăm $\alpha_k = 1/q_k$. Atunci relația (6.1.1) devine

$$S_{k-1} = \alpha_k S_k + \alpha_k r_k \quad (6.1.2)$$

Se propun mai multe scenarii.

6.1.1. Rată constantă, dobândă simplă

Dobânda fiind simplă, dobânda unitară pe intervalul (s, t) este $d(s, t) = i(t-s)$, unde i este rata anuală a dobânzii. Acum $q_j = 1 + d_j = 1 + i(t_j - t_{j-1})$ pentru $j \geq 2$, $q_1 = 1 + it_1$. Scenariul presupune rate egale, adică $r_1 = \dots = r_n = r$.

Rata se calculează punând condiția ca $S_n = 0$.

Relația (6.1.2) implică, din aproape în aproape $S_{n-1} = r\alpha_n \Rightarrow S_{n-2} = r(\alpha_{n-1} + \alpha_{n-1}\alpha_n) \Rightarrow S_{n-3} = r(\alpha_{n-2} + \alpha_{n-2}\alpha_{n-1} + \alpha_{n-2}\alpha_{n-1}\alpha_n) \Rightarrow \dots$ de unde găsim rata

$$r = \frac{S_0}{\alpha_1 + \alpha_1\alpha_2 + \dots + \alpha_1\alpha_2\dots\alpha_n} \quad (6.1.3)$$

Scriptul următor calculează rata r după formula (6.1.3) și amortismentele $(a_k)_{1 \leq k \leq n}$. Atenție: comentariile scrise cu font Times New Roman nu fac parte din script, ci sunt comentarii.

Se apelează prin instrucțiunea

```
a = ratadosi(s0,i,t)
```

Variabilele funcției sunt suma împrumutată s_0 , rata anuală a dobânzii, i și vectorul t cu momentele când se fac plățile.

Rezultatul se dă în forma unei liste. Așadar rata este $a[[1]]$ iar amortismentele sunt $a[[2]]$.

```
#functia ratadosi calculeaza rata constanta daca se da
dobinda
#simpla cu rata i, suma de platit
#si daca se da t=(t1,...,tn) momentele de plata a ratelor
ratadosi<-function(s0,i,t)
{ n=length(t); q=t; alfa=t; bet=t; s=t; amo=t # inițializări
q[1]=1+i*t[1]; for (k in 2:n) {q[k]=1+i*(t[k]-t[k-1])}
##  $q_1 = 1 + it_1, k \geq 2 \Rightarrow q_k = 1 + i(t_k - t_{k-1})$ 
alfa=1/q; for (k in 1:n) {bet[k]=prod(alfa[1:k])}
##  $\alpha_k = 1/q_k, \beta_k = \alpha_1 \dots \alpha_k$ 
r=s0/sum(bet) ##  $r = s_0/(\beta_1 + \beta_2 + \dots + \beta_n)$ 
s[1]=s0*q[1]-r; for (k in 2:n) {s[k]=s[k-1]*q[k]-r} ## Relația (1)
amo[1]=s0-s[1]; for (k in 2:n) {amo[k]=s[k-1]-s[k]} ##  $a_k = S_{k-1} - S_k$ 
rata=list(r,amo) ##  $r = \text{rata}[[1]], \text{amortismente} = \text{rata}[[2]]$ 
rata}
```

Exemplu de aplicare

Se face un credit de 1200 de lei, plătit în 12 rate lunare egale, în regim de dobândă simplă cu termen de grație de 1 an și dobândă 12%. Care este rata și cum evoluează amortismentele?

Soluție. Avem $s_0 = 1200$, $i = 0.12$ și $t = (k/12)_{13 \leq k \leq 24}$. Deci succesiunea de instrucțiuni este

```
>t=13:24/12
> a=ratadosi(1200,.12,t)
> rata=a[[1]];rata
[1] 119.2861
> amortismente=a[[2]];amortismente
[1] -36.7139 106.9190 107.9881 109.0680 110.1587 111.2603
112.3729 113.4966
[9] 114.6316 115.7779 116.9357 118.1050
```

De remarcat că primul amortisment este negativ (s-a acumulat dobânda!) , celelalte cresc în progresie geometrică.

6.1.2. Rată constantă, dobândă compusă

Singura deosebire este că acum $q_1 = (1+i)^{t_1}, k > 1 \Rightarrow q_k = (1+i)^{t_k - t_{k-1}}$,

```
#functia ratadoco calculeaza rata constanta daca se da
#dobinda compusa cu rata i, suma de platit
#si daca se da t=(t1,...,tn) momentele de plata a ratelor
ratadoco<-function(s0,i,t)
{ n=length(t); q=t; alfa=t; bet=t; s=t; amo=t
```



```

q[1]=(1+i)^t[1];for (k in 2:n) {q[k]=(1+i)^(t[k]-t[k-1])}
##  $q_k = (1+i)^{t_k - t_{k-1}}$ ,
alfa=1/q;for (k in 1:n) {bet[k]=prod(alfa[1:k])}
r=s0/sum(bet)
s[1]=s0*q[1]-r;for (k in 2:n) {s[k]=s[k-1]*q[k]-r}
amo[1]=s0-s[1];for (k in 2:n) {amo[k]=s[k-1]-s[k]}
rata=list(r,amo)
rata}

```

Exemplu de aplicare

Se face un credit de 1200 de lei, plătit în 12 rate lunare egale, în regim de dobândă compusă cu termen de grație de 1 an și dobândă 12%. Care este rata și cum evoluează amortismentele?

Soluție.

```

> a=ratadoco(1200,.12,t);a
[[1]]
[1] 119.0274

[[2]]
[1] -37.72552 107.28290 108.30088 109.32853 110.36592
111.41316 112.47034
[8] 113.53755 114.61488 115.70244 116.80031 117.90861

```

Rata este puțin mai mică, deoarece dacă s-ar calcula corect rata lunară a dobânzii, ea nu ar fi egală cu $i/12$, ci cu $(1+i)^{\frac{1}{12}} - 1$. Pentru $i = 0,12$ ar rezulta că rata lunară ar trebui să fie 0.95% și nu 1%, așa cum se ia în mod real.⁸

6.1.3. Amortismente constante, dobândă simplă

Relația folosită este $r_k = a + d_k S_{k-1}$ cu amortismentul $a = S_0/n$. Dobânzile unitare sunt $d_1 = it_1$ și $k \geq 2 \Rightarrow d_k = i(t_k - t_{k-1})$

```

#functia amordosi calculeaza ratele in ipoteza
#amortismentelor constante
# daca se da dobinda simpla cu rata i, suma de platit s0
#si daca se da t=(t1,...,tn) momentele de plata a ratelor
amordosi<-function(s0,i,t)
{n=length(t);d=t;r=t;amo=s0/n
d[1]=i*t[1];for (k in 2:n) {d[k]=i*(t[k]-t[k-1])}
for (k in 1:n) {r[k]=amo+d[k]*s0*(1 - (k-1)/n)}
#  $d_k = i(t_k - t_{k-1})$ 
r}

```

Exemplu de aplicare

Se face un credit de 1200 de lei, plătit în 12 rate lunare, în regim de dobândă simplă cu termen de grație de 1 an și dobândă 12%. Care sunt ratele dacă se planifică 12 amortismente egale a câte 1200 lei?

⁸ Ca o curiozitate, rata zilnică a dobânzii se calculează de unele bănci după formula $i/360$. Ca și cum anul ar avea 360 de zile.

Soluție.

```
t=13:24/12
> a=amordosi(1200,.12,t);a
[1] 256 111 110 109 108 107 106 105 104 103 102 101
```

Cu excepția primeia, ratele formează o progresie aritmetică descrescătoare.

6.1.4. Amortismente constante, dobândă compusă

Singura diferență față de scenariul anterior este că dobânzile unitare sunt acum

$$d_1 = (1+i)^t - 1 \text{ și } k \geq 2 \Rightarrow (1+i)^{t_k - t_{k-1}} - 1.$$

```
#functia amordoco calculeaza ratele in ipoteza
#amortismentelor constante
# daca se da dobinda compusa de rata i, suma de platit s0
#si daca se da t=(t1,...,tn) momentele de plata a ratelor
amordoco<-function(s0,i,t)
{n=length(t);d=t;r=t;amo=s0/n
d[1]=(1+i)^t[1]-1;for (k in 2:n) {d[k]=(1+i)^(t[k]-t[k-1])-1}
for (k in 1:n) {r[k]=amo+d[k]*s0*(1 - (k-1)/n)}
r}
```

Exemplu de aplicare

Se face un credit de 1200 de lei, plătit în 12 rate lunare, în regim de dobândă compusă cu termen de grație de 1 an și dobândă 12%. Care sunt ratele dacă se planifică 12 amortismente egale a câte 1200 lei?

Soluție.

```
> a=amordoco(1200,.12,t);a
[1] 256.7529 110.4377 109.4888 108.5399 107.5910 106.6422
105.6933 104.7444
[9] 103.7955 102.8466 101.8978 100.9489
```

6.1.5. Compararea celor patru scenarii

Dacă dorim să comparăm cele patru scenarii din punctul de vedere al sumei totale pe care o avem de plată, putem apela funcția *comprate*.

```
#functia comprate compara cele patru scenarii prin suma
totala
comprate<-function(s0,i,t)
{n=length(t)
a1=ratadosi(s0,i,t);a1=a1[[1]]
a2=ratadoco(s0,i,t);a2=a2[[1]]
a3=amordosi(s0,i,t)
a4=amordoco(s0,i,t)
S1=n*a1;S2=n*a2;S3=sum(a3);S4=sum(a4)
scenarii<-c("rate egale ds","rate egale dc","amo egale
ds","amo egale dc")
sume<-c(S1,S2,S3,S4)
```

```
rez=data.frame(scenarii,sume)
rez}
```

Exemplu de aplicare

În cazul nostru suma cea mai mică este, după cum rezultă din tabelul de mai jos, de 1419.38 lei, în regim de dobândă compusă, dacă s-ar accepta calculul normal al dobânzii.

```
t=13:24/12; comprate(1200,0.12,t)
              scenarii      sume
1 rate egale dob simpla  1431.433
2 rate egale dob compusa 1428.329
3 amo egale dob simpla   1422.000
4 amo egale dob compusa  1419.379
```

6.1.6. Rată fixă de valoare determinată

Revenim la relația de recurență a sumelor restante, $S_k = q_k S_{k-1} - r_k$ și punem altă întrebare: dacă dorim să plătim o rată fixă, r , câte rate avem de plătit? Ultima rată va fi eventual mai mică. Deci problema este de a calcula numărul $N = \min\{n \mid S_n \leq 0\}$. De exemplu, în scenariu de dobândă compusă, $i = 12\%$, cu termen de gratie de 1 an, câte rate de 100 de lei ne trebuie pentru a achita creditul de 1200 lei?

Presupunem că plățile se fac la momentele $t_0 + nh$. Este posibil ca rata să fie prea mică pentru a putea achita rambursa creditul.

Dacă $q_1 = (1+i)^{t_1} - 1$ și $k \geq 2 \Rightarrow q_k = (1+i)^{t_k - t_{k-1}} - 1 = (1+i)^h - 1 = q$, avem

$S_1 = S_0 q_1 - r$, $S_2 = S_1 q - r = (S_0 q_1 - r)q - r = S_0 q_1 q - r(1+q)$, $S_3 = S_0 q_1 q^2 - r(1+q+q^2)$, ..., și, în general, $S_{n+1} = S_0 q_1 q^n - r(1+q+\dots+q^n)$.

Atunci $S_{n+1} \leq 0 \Leftrightarrow S_0 q_1 q^n \leq r(1+q+\dots+q^n) \Leftrightarrow r(1+\alpha+\dots+\alpha^n) \geq S_0 q_1$, cu $\alpha = 1/q$.

Creditul se poate rambursa într-un număr finit de rate egale cu r dacă există un n ca $S_{n+1} \leq 0$. Trebuie deci ca $r > \frac{S_0 q_1}{1+\alpha+\alpha^2+\dots}$, adică deducem condiția

$$r > S_0 q_1 (1-\alpha) \quad (6.1.4)$$

Funcția care calculează numărul de rate și arată și evoluția sumelor restante este nrates. Funcția e calculată în ipoteza dobânzii compuse, dar cititorul o poate modifica ușor înlocuind q_1 și q . De asemenea, am pus o limită numărului de rate (până la 1000), care poate fi și ea modificată. Sumele restante sunt puse în vectorul s , inițializat cu zerouri. Dacă rata este prea mică, funcția prevede un mesaj în care se dă rata minimă cu care se poate acoperi creditul.

```
#Cite rate in valoare de "r"UM am de platit daca dobinda
este
#compusa, plata se incepe la t0 si se face la intervale de
h ani.
nrrates<-function(r,s0,t0,h,i)
{s=c(rep(0,1000));q1=(1+i)^t0;q=(1+i)^h;alfa=1/q
if(r<=q1*s0*(1-alfa)) {print("rata prea mica,
n=infinite");return()}
s[1]=s0*q1-r;rest=s[1];k=2
while((rest>0)&&(k<1000)){s[k]=s[k-1]*q-r;rest=s[k];k=k+1}}
```

```
rez=c(k-2,s[k-2])
rez=list(rez,s[1:k-1])
rez}
```

Exemplu de aplicare

În cazul nostru avem $r = 100$, $S_0 = 1200$, $t_0 = 1$, $h = 1/12$, $i = 0.12$

```
a=nrrate(100,1200,1,1/12,0.12);a
[[1]]
[1] 14.00000 29.85075

[[2]]
[1] 1244.00000 1155.80406 1066.77124 976.89362 886.16316
794.57178
[7] 702.11130 608.77349 514.55002 419.43248 323.41238
226.48118
[13] 128.63021 29.85075 -69.86600
```

Interpretare: este nevoie de 14 rate, ultima incompletă (de 29.85 lei). După prima rată, suma restantă este de 1244 lei, etc; dacă am plăti 14 rate a 100 lei, ar trebui să primim înapoi 69,87 lei – acesta este sensul ultimului număr din tabelul de mai sus.

Dacă dorim să plătim în rate de 10 lei scriem

```
> a=nrrate(10,1200,1,1/12,0.12);a
```

și obținem reacția

```
[1] "rata prea mica, n=infinit"
[1] 12.63307
```

Ultimul număr ne spune că rata minimă este de 12.64 lei. Dacă am încerca cu rate de 13 lei am obține

```
a=nrrate(13,1200,1,1/12,0.12);a[[1]]
[1] 377.000000 9.707553
```

E nevoie de 33 de ani!

6.1.7. Rată fixă de valoare determinată, dobândă compusă aleatoare

Funcția `nralrate(r,s0,i0,i1)` face o simulare în ipoteza că rata este r , suma împrumutată este s_0 iar rata dobânzii este aleatoare. Mai precis, aici se presupune că ratele se plătesc la intervale constante de timp $\Delta t = h$ și pe fiecare din aceste momente de timp dobânzile unitare formează un $\{i_t\}$ de variabile aleatoare independente și identic repartizate uniform pe intervalul (i_0, i_1) . Pentru a nu risca o buclă infinită, am ales un prag ($n = 1000$) peste care nu mai verificăm evoluția sumelor restante. Acum $N = \min\{n \mid S_n \leq 0\}$ devine o variabilă aleatoare.

```
#presupun ca rata dobinzii este aleatoare. Atunci
#S[n]=S[n-1]qn-r
#dar q[n] sunt presupuse i.i.d. repartizate uniform pe
#intervalul i1,i2
#functia nralrate are ca argumente r (rata), s0 (suma
#imprumutata)
```

```
#i0 (rata minima a dobinzii)
#i1(rata maxima). qn = 1 + i[n] sunt presupusi i.i.d
nralrate<-function(r,s0,i0,i1)
{q1=1+(i1-i0)*runif(1);s[1]=s0*q1-r;rest=s[1];k=2
while((rest>0)&&(k<1000)){q=1+(i1-i0)*runif(1)
s[k]=s[k-1]*q-r;rest=s[k];k=k+1}
rez=c(k-2,s[k-2])
rez=list(rez,s[1:k-1])
rez
```

Exemplu de aplicare

Ca să fim consecvenți, luăm $r = 100$, $S_0 = 1200$, $i_0 = 0.05$, $i_1 = 0.15$.

```
a=nralrate(100,1200,.05,.15);a
[[1]]
[1] 19.000000 7.846417

[[2]]
[1] 1164.531376 1139.006774 1072.575349 1001.426395 993.777986
963.953179
[7] 917.006879 848.128914 749.050657 714.409097 673.291059
585.682507
[13] 504.665078 424.788245 355.531910 284.929749 199.222488
106.591050
[19] 7.846417 -91.942226
```

Interpretare: trebuie 19 rate, din care 18 complete. Vectorul `a[[2]]` dă evoluția sumelor restante.

Putem face o **analiză statistică** a numărului necesar de rate pentru rambursarea creditului cu ajutorul funcției `nralratesim(n,r,s0,i0,i1)`, care prezintă un vector cu n simulări

```
##rata aleatoare, uniform repartizata pe (i0,i1), n simulari
nralratesim<-function(n,r,s0,i0,i1)
{nr=c(rep(0,n))
for (k in 1:n) {nr[k]=nralrate(r,s0,i0,i1)[[1]][1]}
nr}
```

A se remarca rîndul al treilea: funcția `nralrate(...)` returnează o listă cu două elemente: primul este un vector cu două componente (n , S_{n-1}) cu n numărul de rate și S_{n-1} ultima sumă restantă iar a doua componentă este șirul sumelor restante. Dacă lista este a , atunci `a[[1]]` este primul vector iar `a[[1]][1]` este numărul n .

Exemplu:

```
nr=nralratesim(10,100,1200,0.05,.15);nr
[1] 16 18 24 20 16 15 20 19 18 16
```

Dacă dorim să estimăm repartiția variabilei aleatoare N , putem folosi funcția „table”. Iată rezultatul a 10.000 de simulări:

```
> nr=nralratesim(10000,100,1200,0.05,.15)
> a=table(nr);x=as.numeric(names(a));m=length(x);p=x;for (k
in 1:m) {p[k]=a[[k]]}
> x
[1] 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 36
```

```

> p
[1] 4 89 479 1165 1809 2106 1721 1128 731 387 200
97 48 18 10
[16] 6 1 1
> pp=p;for (k in 1:m) {pp[k]=sum(p[1:k])/sum(p)}
> b=rbind(x,pp);b

x 13.0000 14.0000 15.0000 16.0000 17.0000 18.0000 19.0000
20.0000 21.0000
pp 0.0004 0.0093 0.0572 0.1737 0.3546 0.5652 0.7373
0.8501 0.9232

x 22.0000 23.0000 24.0000 25.0000 26.0000 27.0000 28.0000
29.0000 36
pp 0.9619 0.9819 0.9916 0.9964 0.9982 0.9992 0.9998
0.9999 1

```

De exemplu, probabilitatea ca $N \leq 20$ este cam 85%

6.1.8. Probabilitatea ca n rate în valoare de r UM să fie suficiente

Revenim la relația de recurență a sumelor restante, $S_k = q_k S_{k-1} - r_k$. Presupunem, la fel ca în paragraful anterior, că ratele sunt egale între ele, $r_k = r$. Ne punem întrebarea : care este probabilitatea ca n asemenea rate să fie suficiente pentru rambursarea creditului?

Să precizăm.

Lucrurile nu sunt lămurite complet din punct de vedere matematic nici măcar în cazul cel mai simplu. Acesta este următorul: presupunem că ratele anuale ale dobânzilor formează un șir de variabile aleatoare independente și identic repartizate $(i_n)_n$. Atunci factorii de fructificare q_n vor forma și ei un șir de variabile aleatoare independente și identic repartizate, și la fel și coeficienții de actualizare $\alpha_n = 1/q_n$.

O rată este suficientă dacă $S_1 \leq 0 \Leftrightarrow r \geq q_1 S_0 \Leftrightarrow \alpha_1 \geq S_0/r$

Două rate sunt suficiente dacă $S_2 \leq 0 \Leftrightarrow S_0 q_1 q_2 - r(1+q_2) \leq 0 \Leftrightarrow \alpha_1 + \alpha_1 \alpha_2 \geq S_0/r$

Trei rate sunt suficiente dacă $S_3 \leq 0 \Leftrightarrow \alpha_1 + \alpha_1 \alpha_2 + \alpha_1 \alpha_2 \alpha_3 \geq S_0/r$

În general, n rate sunt suficiente dacă

$$S_n \leq 0 \Leftrightarrow \alpha_1 + \alpha_1 \alpha_2 + \dots + \alpha_1 \alpha_2 \dots \alpha_n \geq S_0/r \quad (6.1.5)$$

Probabilitatea ca așa ceva să se întâmple este $F_n(S_0/r)$, unde F_n este funcția de repartiție a variabilei aleatoare $X_n = \alpha_1 + \alpha_1 \alpha_2 + \dots + \alpha_1 \alpha_2 \dots \alpha_n$.

Fie $F = F_1$ funcția de repartiție a coeficientului de actualizare α_1 .

Observăm că $X_2 = X_1(1 + \alpha_2)$. Dacă notăm $Y_2 = 1 + \alpha_2$, atunci Y_2 este independentă de X_1 , are funcția de repartiție $G(x) = F(x-1)$ și $X_2 = X_1 Y_2$.

Cum găsim repartiția lui X_2 ?

Dacă X și Y sunt două variabile aleatoare independente, cu repartițiile F și G , atunci XY are o repartiție H care se obține din F și G printr-o operație numită **convoluție multiplicativă**. Să o notăm cu „ \bullet ”. Deci $H = F \bullet G$. Dacă F și G sunt repartiții discrete, atunci H se poate calcula astfel

$$F \sim \begin{pmatrix} x_1 & x_2 & \dots & x_m \\ p_1 & p_2 & \dots & p_m \end{pmatrix}, G \sim \begin{pmatrix} y_1 & y_2 & \dots & y_n \\ q_1 & q_2 & \dots & q_n \end{pmatrix} \Rightarrow F \bullet G \sim \begin{pmatrix} x_1 y_1 & x_1 y_2 & \dots & x_m y_n \\ p_1 q_1 & p_1 q_2 & \dots & p_m q_n \end{pmatrix}$$

Scriptul convm face exact acest lucru.

```
##Calculeaza convolutia multiplicativa intre (x,p) si
(y,q)
convm<-function(x,p,y,q)
{ xmy=outer(x,y,"*"); pmq=outer(p,q,"*")
a=tapply(pmq,xmy,sum)
x1=as.numeric(names(a))
p1=x1;m=length(x1)
for (i in 1:m) {p1[i]=a[[i]]}
a=list(x1=x1,p1=p1)
a}
```

Variabilele sunt vectorii $x=(x_i)_{1 \leq i \leq m}$, $p=(p_i)_{1 \leq i \leq m}$, $y=(y_j)_{1 \leq j \leq n}$, $q=(q_j)_{1 \leq j \leq n}$.

Scriptul nu verifică dacă p și q sunt vectori de probabilitate, adică dacă sumele sunt 1. Cititorul este invitat să îl modifice în așa fel încât să se dea un mesaj de eroare în caz că p sau q nu sunt de sumă 1.

În primul rând se calculează matricile $xmy_{i,j} = x_i y_j$ și $pmq_{i,j} = p_i q_j$. Aici se vede forța operației `outer` din „R”. În rîndul următor se folosește o instrucțiune și mai puternică, `tapply(pmq, xmy, sum)`. Ea face următorul lucru: aplică matricii `pmq` funcția `sum` depinzînd de `xmy`. Explicit, al doilea argument (adică `xmy`) este văzut ca o colecție de etichete. Se face o listă cu aceste etichete și se adună (funcția `sum`) elementele din `pmq` care au această etichetă.

De exemplu, dacă $x = (-1, 0, 1, 2)$, $y = -x = (1, 0, -1, -2)$, $p = q = (\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$, avem

```
> xmy=outer(x,y,"*"); xmy
      [,1] [,2] [,3] [,4]
[1,]    -1     0     1     2
[2,]     0     0     0     0
[3,]     1     0    -1    -2
[4,]     2     0    -2    -4
```

Deci $xmy_{i,j} = x_i y_j$. Apar 7 valori: -4,-2,-1,0,1,2

```
> pmq=outer(p,q,"*"); pmq
      [,1] [,2] [,3] [,4]
[1,] 0.0625 0.0625 0.0625 0.0625
[2,] 0.0625 0.0625 0.0625 0.0625
[3,] 0.0625 0.0625 0.0625 0.0625
[4,] 0.0625 0.0625 0.0625 0.0625
> a=tapply(pmq,xmy,sum); a
      -4     -2     -1      0      1      2
0.0625 0.1250 0.1250 0.4375 0.1250 0.1250
```

Cea mai frecventă valoare a fost 0: apare în 7 locuri. Deci probabilitatea sa este $7/16 = 0.4375$.

Apelarea funcției se face prin comanda

```
a =convm<-function(x,p,y,q)
```

și rezultatul este o listă, `a`. Primul element este format din vectorul cu elementele $x_i y_j$, sortat în ordine crescătoare iar al doilea cu probabilitățile $p \bullet q$

Exemplu

```
> a=convm(x,p,y,q); a
```

```
[[1]]
[1] -4 -2 -1 0 1 2
```

```
[[2]]
[1] 0.0625 0.1250 0.1250 0.4375 0.1250 0.1250
```

Să presupunem acum că rata dobânzii este $i \sim \begin{pmatrix} 5 & 10 & 15 & 20 \\ .2 & .3 & .3 & .2 \end{pmatrix}$. Ratele se văd în primul rând, exprimate în procente iar probabilitățile lor în al doilea rând. Atunci $q_n \sim \begin{pmatrix} 1.05 & 1.1 & 1.15 & 1.2 \\ .2 & .3 & .3 & .2 \end{pmatrix}$ iar coeficienții de actualizare $\alpha_n \sim \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1.05 & 1.1 & 1.15 & 1.2 \\ .2 & .3 & .3 & .2 \end{pmatrix}$

Vrem să găsim repartiția lui $\alpha_1 + \alpha_1\alpha_2 = F \bullet G$ unde $F = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1.05 & 1.1 & 1.15 & 1.2 \\ .2 & .3 & .3 & .2 \end{pmatrix}$ și $G = \begin{pmatrix} 1 + \frac{1}{1.05} & 1 + \frac{1}{1.1} & 1 + \frac{1}{1.15} & 1 + \frac{1}{1.2} \\ .2 & .3 & .3 & .2 \end{pmatrix}$. Atunci instrucțiunile sunt

```
x=c(1/1.05,1/1.1,1/1.15,1/1.2);y=x+1;p=c(2,3,3,2)/10;q=p
convm(x,p,x+1,p)
$xl
[1] 1.527778 1.557971 1.590909 1.594203 1.625709 1.626984
1.660079 1.666667 1.697723 1.699605 1.735537
[12] 1.746032 1.774892 1.780538 1.818182 1.859410

$pl
[1] 0.04 0.06 0.06 0.06 0.09 0.04 0.09 0.06 0.06 0.09 0.09
0.04 0.06 0.06 0.06 0.04
```

Sau, scris și mai explicit: $(x,p) \bullet (y,q) = (z,r)$ cu z un vector cu 16 valori și r de asemenea cu 16 valori. S-a întâmplat că toate produsele $x_i y_j$ au fost diferite.

Probabilitatea $P(X \geq x)$ este calculată de scriptul tailfrep

```
##coada unei repartitii discrete (x,p)calculata in t
tailfrep<-function(t,x,p)
{tailfrep=1-sum(p[x<t]);tailfrep}
```

De exemplu, dacă dorim să știm care este probabilitatea ca o sumă $S_0 = 1200$ lei să poată fi achitată în două rate $r = 700$ lei, calculăm

```
a=convm(x,p,y,q);prob=tailfrep(1200/700,a[[1]],a[[2]]);prob
```

și răspunsul este prob = 0.35. Dacă mărim rata la 800,

```
a=convm(x,p,y,q);prob=tailfrep(1200/800,a[[1]],a[[2]]);prob
prob = 1.
```

Era și de așteptat.

Cum calculăm celelalte probabilități?

De exemplu $X_3 = \alpha_1 + \alpha_1\alpha_2 + \alpha_1\alpha_2\alpha_3 = \alpha_1(1 + \alpha_2(1 + \alpha_3))$ Observăm că, deoarece variabilele α_i au fost presupuse independente și identic repartizate, $\alpha_2(1 + \alpha_3)$ are

aceeași repartiție ca și X_2 iar α_1 are aceeași repartiție ca și α_3 . Ca atare X_3 are aceeași repartiție ca $\alpha_3(1 + X_2)$.

Mai general, $X_n = \alpha_1(1 + (\alpha_2 + \alpha_2\alpha_3 + \dots + \alpha_2\dots\alpha_n))$ are aceeași repartiție cu $\alpha_n(1 + \alpha_1 + \dots + \alpha_1\dots\alpha_{n-1})$, adică cu $\alpha_n(1 + X_{n-1})$. Dacă notăm faptul că două variabile aleatoare au aceeași repartiție cu $X \stackrel{D}{=} Y$, putem scrie

$$X_n \stackrel{D}{=} \alpha_n(1 + X_{n-1}) \quad (6.1.6)$$

Sau, în termeni de repartiții,

$$F_n = F \bullet h(F_{n-1}) \quad (6.1.7)$$

unde $h(F_{n-1})$ este repartiția variabilei aleatoare $X_{n-1} + 1$.

Scriptul `convmrep` calculează repartiția lui X_n

```
##convolutia multiplicativa pentru alfa: F_n+1 =
F_n*delta(1)convmultF
convmrep<-function(n,x,p)
{y=1;q=1
for (i in 1:n) {a=convm(x,p,y,q);y=a[[1]]+1;q=a[[2]]}
a}
```

De exemplu, pentru a afla repartiția lui X_4 scriem

```
a=convrep(4,x,p);x1=a[[1]];p1=a[[2]]
```

Vectorii `x1` și `p1` vor avea lungimea de 256. Valoarea minimă pentru `x1` este 2.588735 iar cea maximă este 3.545951. Probabilitatea ca 4 rate în valoare de 400 de lei să fie suficiente este dată de comanda

```
a=convmrep(4,x,p);x1=a[[1]];p1=a[[2]];prob=tailfrep(1200/400,x
1,p1);prob
```

și este egală cu 0.5361

Dacă dorim să mărim probabilitatea, mărim puțin rata la 450 lei

```
a=convmrep(4,x,p);x1=a[[1]];p1=a[[2]];prob=tailfrep(1200/450,x1,p1);
prob
[1] 0.9836
```

Deja putem paria că este suficient.

Problema este că nu putem aplica acest algoritm de prea multe ori. Deja la 4 convoluții, șirul de valori pe care îl poate lua X_4 este de 256; în principiu este de 4^n valori. De exemplu, la $n = 10$

```
a=convrep(10,x,p);x1=a[[1]];p1=a[[2]]
```

vectorul `x1` are $n = 1048576$ componente. Calculul durează 5 minute. Nu o să afișăm toate valorile, ci doar un sumar prin instrucțiunea `summary`

```
summary(x1)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
4.192   5.251   5.574   5.599   5.921   7.722
```

Deci minimul este 4.192 iar maximul este 7.722.

Probabilitatea ca acest credit să fie acoperit în 10 rate a 200 de lei este

```
prob=tailfrep(1200/200,x1,p1);prob = 0.178917
```

iar în rate de 250 lei este

```
prob=tailfrep(1200/250,x1,p1);prob = 0.9737461
```

Funcția $\text{gfrep}(n,x,p)$ face graficul funcției de repartiție a repartiției (x,p) în n puncte echidistante calculate între $\min x_1$ și $\max x_1$. Iată graficul funcției de repartiție al lui X_{10}

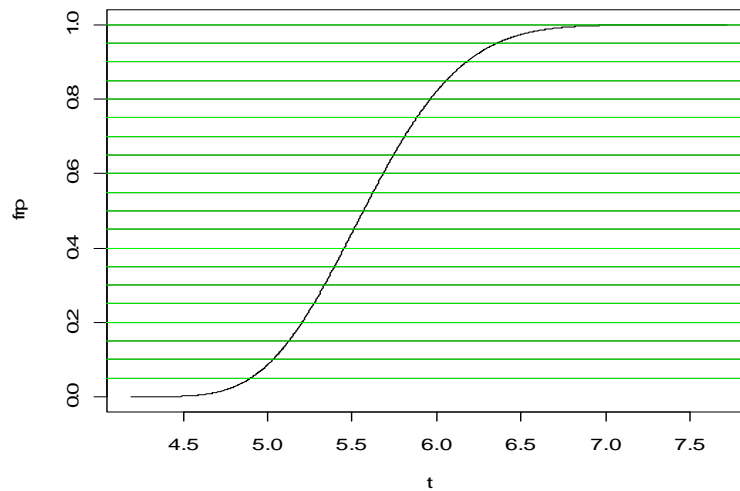


Fig 6.1

Pe axa Ox sunt valorile lui x , între minim și maxim. În graficul următor prezentăm diferențele $[F(t+h) - F(t)]/h$ pentru $h = (\max x_1 - \min x_1)/998$ care, în ipoteza că repartiția lui F_n ar avea o limită absolut continuă, ar mima densitatea sa. Dar nu știm dacă există o asemenea limită: este o problemă nerezolvată de matematicieni

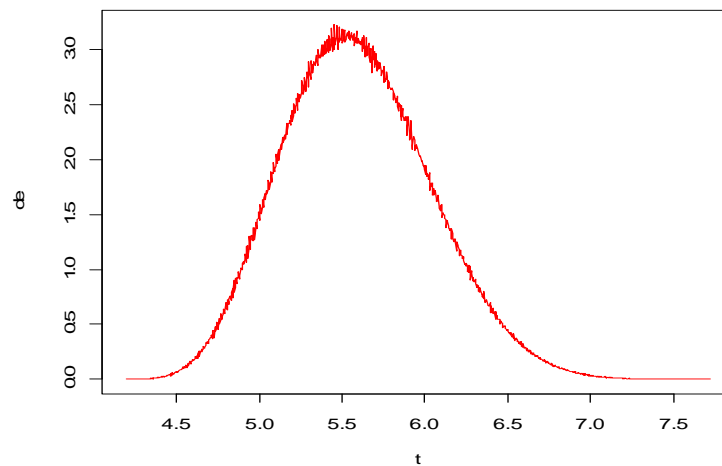


Fig 6.2

Pentru valori mai mari ale lui n nu este realist să aplicăm acest algoritm: dacă la $n = 10$ X_n are 1048576 de componente, la $n = 20$ o să aibă 10^{12} . Am putea mări numărul de iterații dacă am modifica problema: dacă rata dobânzii nu are decât două valori $i \sim \begin{pmatrix} 0.10 & 0.15 \\ .5 & .5 \end{pmatrix}$. Ratele se văd în primul rând, exprimate în procente iar probabilitățile lor în al doilea rând. Atunci $q_n \sim \begin{pmatrix} 1.10 & 1.15 \\ .5 & .5 \end{pmatrix}$ iar coeficienții de actualizare $\alpha_n \sim \begin{pmatrix} 1 & 1 \\ 1.10 & 1.15 \\ .5 & .5 \end{pmatrix}$. Atunci am putea calcula repartiția chiar și pentru X_{20}

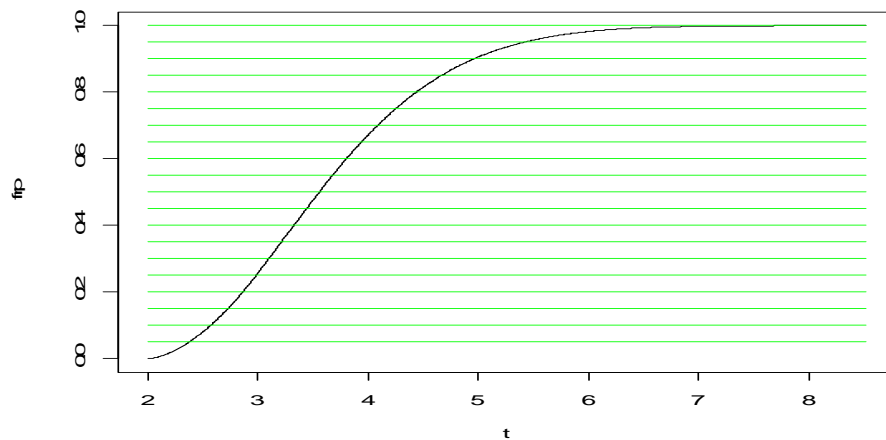


Fig 6.3

și „densitatea”

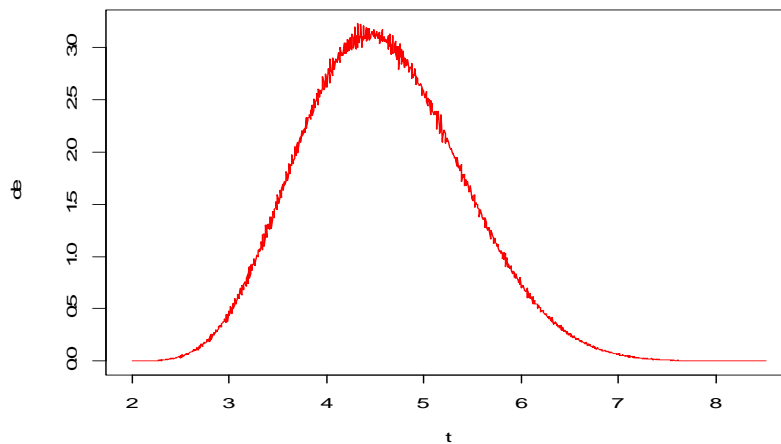


Fig 6.4

Probabilitatea ca 20 de rate a 100 de lei să fie suficiente este

```
prob=tailfrep(1200/250,x1,p1);prob = 0.1249514
```

6.1 9. Calculul mediei și varianței variabilelor $X_n = \alpha_1 + \alpha_1 \alpha_2 + \dots + \alpha_1 \alpha_2 \dots \alpha_n$

Dacă n este mare – de ordinul zecilor, atunci algoritmul de mai sus nu poate fi aplicat dacă ținem ca modelul să fie cât de cât realist. Dar dacă presupunem că variabilele $\alpha_i \in (0,1)$ sunt independente și identic repartizate, atunci putem încerca estimări bazate pe inegalitatea lui Cebîșev. Chiar dacă sunt slabe, ne dau o idee despre repartiție. Să presupunem că

$$E\alpha_i = a, E\alpha_i^2 = b, EX_n = \mu, \text{Var}(X_n) = v \quad (6.1.8)$$

Datorită faptului că variabilele $(\alpha_i)_i$ au fost presupuse independente, avem

$$\mu = E\alpha_1 + E\alpha_1 E\alpha_2 + \dots + E\alpha_1 E\alpha_2 \dots E\alpha_n = a + a^2 + \dots + a^n, \text{ adică}$$

$$\mu = \frac{a(1-a^n)}{1-a} \quad (6.1.9)$$

Pe de altă parte $v = \text{Var} \sum_{i=1}^n \alpha_1 \dots \alpha_i = \sum_{i,j} c_{i,j}$ unde $c_{i,j} = \text{cov}(\alpha_1 \dots \alpha_i, \alpha_1 \dots \alpha_j) = E(\alpha_1 \dots \alpha_i \cdot \alpha_1 \dots \alpha_j) - a^{i+j}$. Dacă $i < j$ avem $c_{i,j} = b^i a^{j-i} - a^{i+j}$, deci

$$c_{i,j} = b^{\min(i,j)} a^{|i-j|} \quad (6.1.10)$$

În concluzie, după efectuarea calculelor ajungem la

$$v = \frac{2a}{1-a} \left(\frac{b(1-b^{n-1})}{1-b} - \frac{ab(a^{n-1}-b^{n-1})}{a-b} \right) + \frac{b(1-b^n)}{1-b} - \left(\frac{a(1-a^n)}{1-a} \right)^2 \quad (6.1.11)$$

Să presupunem, ca la exemplul de la paragraful 7, că rata dobânzii este repartizată Uniform(0.05, 0.15)

Atunci $q \sim \text{Uniform}(1.05, 1.15)$. Cum $U \sim U(\alpha, \beta)$, $\Rightarrow E \frac{1}{U} = \frac{\ln \beta - \ln \alpha}{\beta - \alpha}$ și

$$E \frac{1}{U^2} = \frac{1}{\alpha \beta} \text{ avem, în cazul nostru, } a = \frac{\ln \beta - \ln \alpha}{\beta - \alpha}, b = \frac{1}{\alpha \beta}$$

`alfa=1.05; beta=1.15; a=(log(beta)-log(alfa))/(beta-alfa); b=1/alfa/beta`

Rezultă $a = 0.9097178$, $b = 0.8281573$

Scriptul următor, `medvar(n,a,b)` calculează media μ și varianța v

```
##calculez media si varianta lui X_n
#=#alfa_1+alfa_1alfa_2....
##daca Ealfa_n = a, Ealfa_n^2=b
medvar<-function(n,a,b){
mu=a*(1-a^n)/(1-a)
s=b*(1-b^(n-1))/(1-b)-a*b*(a^(n-1)-b^(n-1))/(a-b)
v=2*a*s/(1-a)+b*(1-b^n)/(1-b)-(a*(1-a^n)/(1-a))^2
answ=list(rate=n,media=mu,varianta=v)
answ}
```

De exemplu, dacă $n = 60$, atunci X_{60} are media 10.04189 și varianța 0.402313

```
ans=medvar(60,a,b);ans
$rate
[1] 60
```

```
$media
[1] 10.04189
```

```
$varianta
[1] 0.402312
```

Atunci abaterea medie pătratică σ este 0.6342815. Deci $P(10.02 - k \cdot 0.63 < X_{20} < 10.02 + k \cdot 0.63) > 1 - 1/k^2$. O aproximație bună este că $P(10.02 - k \cdot 0.63 < X_{20}) > 1 - \frac{1}{2k^2}$. Pentru $k = 4$, găsim că $P(X_{20} > 7.503) > 96.88\%$. Dacă $S_0 = 1200$, punem condiția ca $1200/r < 7.5 \Rightarrow r > 160$. O rată de 160 lei ar fi suficientă, cu probabilitate mai mare decât 97% ca să achite creditul de 1200 lei în 20 de rate.

6.2. Probabilitatea de ruină

O companie plătește în fiecare zi/săptămână/lună (depinde cum se face balanța) X_n lei și câștigă Y_n lei. Pierderea sa zilnică este atunci $\xi_n = X_n - Y_n$. Pierderea după n zile este

$$S_0 = 0, \quad n \geq 1 \Rightarrow S_n = \xi_1 + \dots + \xi_n \quad (6.2.1)$$

Dacă admitem (o simplificare matematică!) că variabilele $(X_n)_n$ sunt i.i.d., variabilele $(Y_n)_n$ sunt și ele i.i.d. și că X_n este independentă de Y_n pentru orice n , atunci repartiția variabilei aleatoare S_n este convoluția repartițiilor variabilelor aleatoare ξ_n . Iar repartiția lui ξ_n este convoluția dintre repartiția F_X a variabilelor X_n și repartiția F_Y a variabilelor $-Y_n$. Convoluția se notează cu „*”.

Vom considera repartițiile F_X și F_Y ca fiind discrete – altfel nu prea putem folosi calculatorul. Atunci $F_X = \begin{pmatrix} x_1 & x_2 & \dots & x_m \\ p_1 & p_2 & \dots & p_m \end{pmatrix}$, $F_Y = \begin{pmatrix} y_1 & y_2 & \dots & y_m \\ p_1 & p_2 & \dots & p_m \end{pmatrix}$ și rezultă că

$$F_X * F_Y = \begin{pmatrix} x_1 + y_1 & x_1 + y_2 & \dots & x_m + y_m \\ p_1 q_1 & p_1 q_2 & \dots & p_m q_n \end{pmatrix}. \text{ Fiind discrete, cele două repartiții pot fi}$$

codificate prin (x, p) și (y, q) . În cazul acesta putem scrie $(x, p) * (y, q) = (x_1, p_1)$

Scriptul `conva(x, p, y, q)` calculează exact acest lucru: produce vectorii x_1 și p_1 . Este o operație mai obișnuită decât convoluția multiplicativă prezentată în capitolul anterior. Diferența este că valorile x_i și y_j se **adună**, nu se înmulțesc.

```
##convolutia aditiva obisnuita cu functia tapply
conva<-function(x, p, y, q)
{ xay=outer(x, y, "+"); pmq=outer(p, q, "*")
a=tapply(pmq, xay, sum)
x1=as.numeric(names(a))
p1=x1;m=length(x1)
for (i in 1:m) {p1[i]=a[[i]]}
a=list(x1=x1, p1=p1)
a}
```

Să presupunem, ca să revenim la scenariul nostru, că $X_n \sim \begin{pmatrix} 0 & 1 & 3 & 10 \\ .3 & .4 & .2 & .1 \end{pmatrix}$.

Probabilitatea să avem o zi fără cheltuieli este 30%, să cheltuim 1 UM este 40%, să cheltuim 3 UM este 20% și, în zilele rele, putem cheltui chiar 10 UM, cu probabilitatea

10%. Câștigurile le presupunem a avea repartiția $Y_n \sim \begin{pmatrix} 1 & 2 & 3 & 4 \\ .4 & .3 & .2 & .1 \end{pmatrix}$. Am ales exemplul în așa fel încât $EX_n = EY_n = 2$. Care este repartiția lui ξ_n ?

Instrucțiunile vor fi

```
x=c(0,1,3,10);y=c(1,2,3,4);p=c(3,4,2,1)/10;q=c(4,3,2,1)/10
a=conva(x,p,-y,q);a
$xl -4 -3 -2 -1 0 1 2 6 7 8 9
$p1 0.03 0.10 0.17 0.26 0.20 0.06 0.08 0.01 0.02 0.03 0.04
```

Probabilitatea să nu avem pierderi este $P(\xi_n \leq 0)$, pentru care avem scriptul

```
frep<-function(t,x,p)
{frep=sum(p[x<=t])frep}
```

Rezultă concret $frep(0,xl,p1)=0.76$

Dacă dorim să calculăm repartiția sumei S_n apelăm funcția convarep. De exemplu, pentru $n = 3$ avem

```
a=conva(x,p,-y,q);xl=a[[1]];p1=a[[2]];a=convarep(3,xl,p1);a
$xn
[1] -12 -11 -10 -9 -8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5
6
[20] 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
25 26 27
$yn
[1] 0.000027 0.000270 0.001359 0.004762 0.012921 0.028188 0.050933 0.077658
[9] 0.100392 0.111242 0.106719 0.088602 0.065603 0.045132 0.030813 0.027390
[17] 0.031509 0.037536 0.041444 0.038796 0.030717 0.020202 0.012105 0.006282
[25] 0.003435 0.003792 0.004395 0.005130 0.004788 0.003618 0.001993 0.000870
[33] 0.000405 0.000056 0.000111 0.000174 0.000219 0.000204 0.000144 0.000064
```

Pentru probabilitatea de a nu fi în pierdere

```
xl=a[[1]];p1=a[[2]];prob=frep(0,xl,p1);prob = 0.648676
```

Observăm că este mai mică.

Să presupunem acum că avem o rezervă de u UM. Dacă pierderea depășește u , declarăm faliment. **Care este probabilitatea de faliment după n zile?**

Putem reformula problema în termeni de pierdere maximă.

Fie $L_n = \max(S_0, S_1, \dots, S_n)$. Aceasta este **pierderea maximă** pe care o avem după n zile. Ruina apare dacă $L_n > u$. Probabilitatea $P(L_n > u)$ se notează cu $\psi_n(u)$ și se numește **probabilitatea de ruină** după n pași. Observăm că ea este coada funcției de repartiție a pierderii maxime L_n .

Fie $F = F_X * F_Y$ repartiția pierderii ξ_n și F_n repartiția sumei S_n .

Observăm că $L_1 = (\xi_1)_+$. Dacă o variabilă aleatoare are repartiția F , notăm cu F_+ repartiția părții sale pozitive. Adică $F_+(x) = F(x)$ dacă $x \geq 0$ și $F_+(x) = 0$ dacă $x < 0$.

Așadar $F_1 = F_+$

Avem apoi $L_2 = \max(0, \xi_1, \xi_1 + \xi_2) = \max(0, \xi_1 + 0, \xi_1 + \xi_2) = \max(0, \xi_1 + \max(0, \xi_2))$. Dar ξ_1 are aceeași repartiție cu ξ_2 iar $\max(0, \xi_2)$ are aceeași repartiție ca și $L_1 = (\xi_1)_+$, deci L_2 are aceeași repartiție ca $(L_1 + \xi_2)_+$. În concluzie, repartiția sa va fi $F_2 = (F_1 * F)_+$.

Analog avem $L_3 = (L_2 + \xi_3)_+$ și, în general rezultă recurența

$$L_{n+1} = (L_n + \xi_{n+1})_+ \Leftrightarrow F_{n+1} = (F_n * F)_+ \quad (6.2.2)$$

Scriptul care calculează partea pozitivă a unei repartiții este ppos

```
## partea pozitiva a unei repartitii
ppos<-function(x,p)
{a=data.frame(v=c(0,x[x>0]),p=c(sum(p[x<=0]),p[x>0]))
a}
```

cel care calculează partea pozitivă a unei convoluții este convap

```
##convolutia pozitiva a doua repartitii
convap<-function(x,p,y,q)
{a=conva(x,p,y,q);x1=a[[1]];p1=a[[2]]
a=ppos(x1,p1)
a}
```

iar cel care calculează convoluția pozitivă iterată, adică exact repartiția F_n din (2.2) este convaprep

```
##convolutia pozitiva iterata pentru Fn+1=(F*Fn)+
# (x,p) este pierderea, (y,p) e cistigul
# se convoluteaza la inceput (x,p) cu (-y,q) si rezulta F
convaprep<-function(n,x,p,y,q)
{y=-y;z=conva(x,p,y,q)
x=z[[1]];p=z[[2]] #aici s-a facut F
y=0;q=1
for (i in 1:n) {z1=convap(x,p,y,q);y=z1[,1];q=z1[,2]}
z1}
```

Astfel, dacă dorim să calculăm repartiția pierderii maxime după 10 zile, folosim comanda

```
a=convaprep(10,x,p,y,q);a
      v      p
0 3.176752e-01
1 6.331282e-02
2 6.736288e-02
3 4.082990e-02
4 4.043586e-02
5 4.000655e-02
6 4.474262e-02
7 4.794617e-02
8 4.961629e-02
9 4.887334e-02
10 3.192275e-02
11 2.745033e-02
12 2.284033e-02
13 2.112670e-02
14 1.962573e-02
15 1.826939e-02
16 1.636947e-02
17 1.408383e-02
18 1.158424e-02
19 9.172632e-03
.....
```

Rezultatul este prezentat sub forma unei matrici $N \times 2$, unde N este numărul de valori pe care le poate lua S_n . În cazul nostru, $N = 91$.

Să presupunem că managerul companiei are la dispoziție un capital $u = 20$ UM pentru a plăti cheltuielile. Care este probabilitatea ca acest capital să nu fie suficient?

Este $\psi_{10}(20) = P(L_{10} > 20)$. Pentru ea avem la dispoziție funcția `tailfrep`

```
prob=tailfrep(20,x1,p1);prob = 0.04675292
```

Dar după 100 de zile?

```
> tailfrep(20,x1,p1); prob =0.4706843
```

Deja probabilitatea este aproape de $\frac{1}{2}$.

După 200 de zile, probabilitatea este $\psi_{200}(20) = 0.6063088$

De curiozitate, $\psi_{200}(100) = 0.0224711$. Probabilitatea de ruina este de doar 2.2% dacă se dispune de o rezervă de 100 lei.

Se poate arăta că dacă $EX \geq EY$ (acesta este cazul nostru), **ruina este sigură**. Adică $\psi_n(u) \rightarrow 1$ dacă $n \rightarrow \infty$ indiferent de capitalul inițial u . Este adevărat că dacă $EX = EY$ se converge lent.

Ruina nu este sigură doar dacă $EX < EY$.

Să înlocuim în exemplul nostru probabilitățile în așa fel încât

$X_n \sim \begin{pmatrix} 0 & 1 & 3 & 10 \\ .6 & .2 & .1 & .1 \end{pmatrix}$. Acum $EX = 1.5$, $E\xi_n = -0.5$, deci ruina nu mai este

sigură. Repartiția lui L_{100} e calculată cu instrucțiunile

```
> x=c(0,1,3,10);p=c(6,2,1,1)/10;a=convaprep(100,x,p,y,q)
> x1=a[,1];p1=a[,2]
> summary(x1)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
    0      225     450     450     675     900
> tailfrep(20,x1,p1); prob = 0.135263
```

Un capital inițial de doar 20 UM face ca probabilitatea de ruină după 100 de zile să fie în jur de 13.5%.

Putem compara această probabilitate cu marginea dată de Lundberg.

Notăm $\psi(u) = \lim_{n \rightarrow \infty} \psi_n(u)$. Marginea lui Lundberg este $\psi(u) \leq e^{-Ru}$ unde R este soluția ecuației $m_\xi(t) = 1$.

```
##calculul functiei generatoare de momente
mgf<- function(t,x,p)
{y=exp(t*x); mgf=y%%p
mgf}
```

Soluția este $R = -0.08551$

Marginea Lundberg: $\psi(20) \leq e^{-20u} = e^{-1.7102} = 0.1808296$

Marginea Lundberg este una pesimista.

6. 3. Probleme de asigurări. Suma daunelor, recurența Panjer

Pentru un asigurator, fiecare client reprezintă o cerere potențială de daună, deci el trebuie să fie pregătit să își onoreze obligațiile. Ce rezervă trebuie să aibă?

În cel mai simplu model, clienții sunt modelați printr-un șir de variabile aleatoare și identic repartizate $(X_n)_n$ cu o funcție de repartiție F . Acesta se numește un **portofoliu omogen de tip F** . În planificarea rezervei, mărimea N a portofoliului (adică numărul de clienți de tip F) este necunoscută. Este o altă variabilă aleatoare N cu valori numere naturale. În modelul cel mai simplu, contorul N va fi independent de variabilele $(X_n)_n$.

Atunci suma pe care trebuie să o aibă la dispoziție managerul va fi o variabilă aleatoare

$$S_N = X_1 + X_2 + \dots + X_N \text{ (dacă } N \geq 1) \text{ sau } S_N = 0 \text{ (dacă } N = 0) \quad (6.3.1)$$

Asiguratorul își va alege un risc α (de exemplu $\alpha = 5\%$, $\alpha = 1\%$ etc) și va aloca portofoliului de tip F o suma Σ în așa fel ca $P(S_N > \Sigma) \leq 1 - \alpha$

Ca să fie posibil acest lucru, va trebui să știe repartiția lui S_N . De regulă, acest lucru este imposibil. O primă abordare ar fi să se folosească inegalitatea lui Cebîșev

$$P(|S_N - ES_N| > k\sigma(S_N)) < 1/k^2 \quad (6.3.2)$$

Sau, același lucru, dar pus în formă afirmativă

$$P(ES_N - k\sigma(S_N) \leq S_N \leq ES_N + k\sigma(S_N)) > 1 - 1/k^2 \quad (6.3.3)$$

pentru că există formule pentru ES_N și $\text{Var}(S_N)$, numite formulele lui Wald ⁹

$$ES_N = EX_1 \cdot EN, \text{Var}(S_N) = \text{Var}(X_1)EN + \text{Var}(N)E^2X_1 \quad (6.3.4)$$

Exemplul 1. Să spunem că $X_n \sim \text{Bin}(10, 1/2)$ și $N \sim \text{Poisson}(100)$. Ne așteptăm să avem în jur de 100 de clienți și fiecare să ceară, în medie 5 lei. Ce sumă să avem pregătită?

Avem $EX_1 = 10 \cdot 1/2 = 5$, $\text{Var}(X_1) = 10 \cdot 1/2 \cdot (1 - 1/2) = 2.5$, $EN = 100$, $\text{Var}(N) = 100$

Răspuns. $ES_N = 5 \cdot 100 = 500$, $\text{Var}(S_N) = 2.5 \cdot 100 + 100 \cdot 5^2 = 2750 \Rightarrow \sigma(S_N) = \sqrt{2750} \approx 52.44$

Aplicînd (6.3.4) vedem că $P(500 - k \cdot 52.44 \leq S_N \leq 500 + k \cdot 52.44) > 1 - 1/k^2$

Alegînd $k = 10$, avem $P(S_N \leq 1024) > 0.99$. Deci ar trebui planificați 1024 de lei.

Inegalitatea lui Cebîșev poate fi îmbunătățită, dacă înlocuim varianța cu semivarianța superioară $\text{Var}_+(X) = E(X - EX)_+^2$ și abaterea standard cu abaterea superioară $\sigma_+(X)$.

S-a constatat empiric faptul că nu se greșește prea mult dacă se ia $\text{Var}_+(X) = \text{Var}(X)/2 \Leftrightarrow \sigma_+ = \sigma/\sqrt{2}$

Inegalitatea (3.3) devine

$$P(S_N \leq ES_N + k\sigma(S_N)/\sqrt{2}) > 1 - 1/k^2 \quad (6.3.5)$$

În cazul exemplului de mai sus, ar rezulta $P(S_N \leq 500 + 370.8) > 0.99 \Leftrightarrow P(S_N \leq 871) > 0.99$

Parcă ar fi suficienți 871 lei, nu 1024. Poate că e nevoie și de mai puțini?

Dacă dorim să calculăm o limită și mai precisă, ar trebui să putem calcula repartiția lui S_N .

Uneori acest lucru este posibil. Dacă acceptăm că X_n au o repartiție discretă,

⁹ Vezi, de exemplu Gheorghița Zbăganu, *Metode matematice în teoria riscului și actuariat*, Editura Universității București, 2004

$X_n \sim \begin{pmatrix} 0 & 1 & \dots & k \\ p(0) & p(1) & \dots & p(k) \end{pmatrix}$ și că $N \sim \begin{pmatrix} 0 & 1 & 2 & \dots \\ q(0) & q(1) & q(2) & \dots \end{pmatrix}$, atunci S_N
 $\sim \begin{pmatrix} 0 & 1 & 2 & \dots \\ f(0) & f(1) & f(2) & \dots \end{pmatrix}$ și probabilitățile $f(n)$ se pot calcula.
 Într-adevăr,

$$f(j) = P(S_N = j) = \sum_{n=1}^j P(S_n = j) q_j \quad (6.3.6)$$

Pentru probabilitățile $P(S_n = j)$ avem la dispoziție scriptul convarep

```

Convarep<-function(n,x,p)
{y=0;q=1
for (i in 1:n){a=conva(x,p,y,q);y=a[[1]];q=a[[2]]}
xn=a[[1]];yn=a[[2]]
a=list(xn=xn,yn=yn)}

```

care calculează a n –a convoluție a repartiției (x,p) a lui X_n .

Dificultatea în programare nu ar fi așa de mare; facem o matrice $pr_{j,i} = P(S_j = i)P(N = i)$ și calculăm sumele de la (3.6)

Iată scriptul convarepN (convoluții aditive repetate de N ori) care calculează probabilitățile f_j și funcția de repartiție $F(j) = f_0 + \dots + f_j$

```

## functia convarepN calculeaza repartitia sumei S_N
## daca nici x nici N nu au prea multe valori
convarepN<-function(x,p,N,q){
nx=length(x);nN=length(N);nmax=nx*nN;f=c(rep(0,nmax));F=f
pr=c(rep(0,nN*nmax));dim(pr) = c(nN,nmax)
pr[1,1]=q[1]
for (j in 2:nN){a=convarep(j-1,x,p);p1=a[[2]];nj=length(p1)
for (i in 1:nj){
pr[j,i]=p1[i]*q[j]
}}
for (j in 1:nmax){f[j]=sum(pr[,j])}
for (j in 1:nmax){F[j]=sum(f[1:j])}
a= list(f=f,F=F)a}

```

Rezultatul este prezentat într-o listă. Apelarea se face prin comanda

```
a= convarepN(x,p,N,q);f=a[[1]];F=a[[2]]
```

Problema este că prin apelarea scriptului convarep se face un consum mare de resurse. De exemplu, dacă $N \sim \text{Uniform}(\{0,1,\dots,40\})$ și $X_n \sim \text{Binom}(10, \frac{1}{2})$, rularea scriptului convarepN durează 2 minute. Scriptul se poate optimiza, renunțând la scriptul convarep, dar programul se mai complică. Putem compara cât de eficientă este metoda inegalității lui Cebîșev față de cea exactă.

Exemplul 2. $X_n \sim \text{Binom}(10, \frac{1}{2})$, $N \sim \text{Uniform}(\{0,1,\dots,40\}) \Rightarrow EX = 5$, $\text{Var}(X) = 2.5$, $EN = 20$, $\text{Var}(N) = 140$, $ES_N = 5 \cdot 20 = 100$, $\text{Var}(S_N) = \text{Var}(X_1)EN + \text{Var}(N)E^2X_1 = 2.5 \cdot 20 + 140 \cdot 5^2 = 3550$, așadar $ES_N = 100$, $\sigma(S_N) \approx 59.58$. Formula (3.5) ne dă $P(S_N \leq 100 + k \cdot 42.13075) > 1 - 1/k^2$. Dacă dorim un interval de predicție cu risc 1% suntem forțați să luăm $k = 10$ și obținem predicția $S_N \leq 521$. Adică managerul de care

era vorba la început să aibă pregătiți 521 de lei pentru a avea doar riscul 1% ca suma să fie insuficientă.

Calculul exact reclamă calculul cuantilei de 0.99 (sau, cum i se mai spune, valoarea la risc, adică primul k cu proprietatea că $F(k) \geq 0.99$

Instrucțiunea care face acest lucru este `which(F<.99)` și obținem 207 lei. Într-adevăr, $F(207) = 0.990797$, $F(206) = 0.9892509$ ¹⁰. De altfel, marginea dată de inegalitatea Cebîșev este absurdă, căci S_N nu ia decât 451 de valori.

Există cazuri în care calculul repartiției lui S_N se poate face cu mult mai rapid, și cu resurse mult mai puține, prin algoritmul lui Panjer.

Definiție. Variabila aleatoare $N \sim \begin{pmatrix} 0 & 1 & 2 & \dots \\ q_0 & q_1 & q_2 & \dots \end{pmatrix}$ se numește **contor de tip**

Panjer dacă $q_0 > 0$ și există $a, b \in \mathfrak{R}$ în așa fel încât $p_n = p_{n-1} \left(a + \frac{b}{n} \right)$ ¹¹

Există trei tipuri de contoare de tip Panjer

- $N \sim \text{Bin}(n, \alpha) \Rightarrow a = -\frac{\alpha}{1-\alpha}, b = \frac{\alpha(n+1)}{1-\alpha}$
- $N \sim \text{Poisson}(\lambda) \Rightarrow a = 0, b = \lambda$
- $N \sim \text{Negbin}(n, \alpha) \Rightarrow a = 1 - \alpha, b = (1 - \alpha)(n-1)$

Panjer a arătat că dacă $X_n \sim \begin{pmatrix} 0 & 1 & \dots & k \\ p_0 & p_1 & \dots & p_k \end{pmatrix}$ atunci numerele $f_n = P(S_N = n)$

satisfac recurența

$$f_0 = g(p_0), n \geq 1 \Rightarrow f_n = \frac{1}{1 - ap_0} \sum_{j=1}^n \left(a + \frac{bj}{n} \right) p_j f_{n-j} \quad (6.3.7)$$

Aici $g(x) = \sum_{n=0}^{\infty} q_n x^n$ este funcția generatoare a contorului N .

Pe baza acestei recurențe se poate face un algoritm de calcul rapid al convoluției F^{*n} pe care îl prezentăm în scriptul `convrap(x, p, N, epsilon)`.

Ideea este că dacă $X_n \sim \begin{pmatrix} 0 & 1 & \dots & k \\ p_0 & p_1 & \dots & p_k \end{pmatrix}$ și $Y_n = (X_n \mid X_n > 0) \sim \begin{pmatrix} 1 & 2 & \dots & k \\ \frac{p_0}{1-p_0} & \frac{p_1}{1-p_0} & \dots & \frac{p_k}{1-p_0} \end{pmatrix}$,

atunci $X_1 + \dots + X_n \stackrel{D}{=} Y_1 + \dots + Y_N$ unde $N \sim \text{Binom}(n, 1 - p_0)$ este independent de (X_n) .

În cuvinte, la adunare nu contează zerourile, iar numărul variabilelor aleatoare X_k diferite de 0 este o variabilă aleatoare $N \sim \text{Binom}(n, P(X_1 > 0)) = \text{Binom}(n, 1 - p_0)$.

¹⁰ Atragem atenția că $F[k]$ nu calculează $F(k)$, ci $F(k-1)$, deoarece în „R” indicii vectorilor încep cu 1, nu cu 0 ca în alte medii de programare (de exemplu Basic).

¹¹ Vezi de exemplu G. Zbăganu, *Metode Matematice în teoria riscului*..., pg. 170

¹² Într-adevăr, dacă $Y = (X \mid X > 0)$, $Ee^{tY} = E(e^{tX} \mid X > 0) = E(e^{tX}; X > 0) / (1 - p_0) = (m_X(t) - p_0) / (1 - p_0)$. Fie $S = X_1 + \dots + X_n$ și $T = Y_1 + \dots + Y_N$. Atunci, cum N este independent de X , $Ee^T = g_N(m_Y) = (p_0 + (1 - p_0)m_Y)^n = \left(p_0 + \frac{m_X - p_0}{1 - p_0}(1 - p_0) \right)^n = (m_X)^n = Ee^S$. Adică $S \stackrel{D}{=} T$.

Deoarece $g_N(x) = (p_0 + (1-p_0)x)^n$, $f_0 = p_0^n$, $a = -\frac{1-p_0}{p_0}$, $b = \frac{(1-p_0)(n+1)}{p_0}$ iar recurența devine $f_k = \frac{1}{1-ap(Y=0)} \sum_{j=1}^k \left(a + \frac{bj}{k}\right) p(Y=j) f_{k-j} = \frac{1}{p_0} \sum_{j=1}^k \left(-1 + \frac{(n+1)j}{k}\right) p_j f_{k-j}$,

deci

$$f_0 = p_0^n, k \geq 1 \Rightarrow f_k = \frac{1}{kp_0} \sum_{j=1}^k ((n+1)j - k) p_j f_{k-j} \quad (6.3.8)$$

Vom face un script care să calculeze f_S și F_S , probabilitățile și funcția de repartiție ale lui S . Apar două probleme de programare:

- în principiu, cu excepția primului caz, S va avea o infinitate de valori. Nici un calculator nu poate calcula așa ceva. Deci trebuie să ne hotărîm pînă unde calculăm probabilitățile f_k . Vom alege un ε și vom calcula f_k pînă la primul k la care $F(k) > 1 - \varepsilon$
- în “R”, vectorii sunt numerotați cu indici începînd de la 1, și nu de la 0. Trebuie mare atenție ca să nu ne încurcăm în indici. Dacă x, p, f sunt gîndiți ca vectori, atunci $f_j = f[j+1]$, $p_j = p[j+1]$ și recurența devine

$$f[1] = g(p[1]), f[n+1] = \frac{1}{1-ap[1]} \sum_{j=1}^n \left(a + \frac{bj}{n}\right) p[j+1] f[n-j+1]$$

- Totuși, f și F trebuie inițializați ca vectori de o anumită lungime. Am ales această lungime să fie egală cu $n_{\max} = ES + 15\sigma(S)$. Inegalitatea lui Cebîșev spune că $P(S > n_{\max}) < 1/5^2 = 0.44\%$. În realitate ea este mult mai mică.

Scriptul următor face acest lucru. El este gîndit ca o funcție Panjer(x,p,tip,N,pN,epsilon).

Mai întîi verificăm dacă datele de intrare sunt corecte: x trebuie să fie un vector cu numere naturale, p un vector cu numere pozitive de sumă 1 de aceeași lungime ca și x iar pN trebuie să fie un număr cuprins între 0 și 1.

Apoi standardizăm pe x și p . Sortăm pe x în ordine crescătoare și completăm la vectorul p locurile libere 0. De exemplu, dacă $x = (5, 3, 6)$, $p = (0.5, 0.2, 0.3)$ scriem $x = (0, 1, 2, 3, 4, 5, 6)$, $p = (0, 0, 0, .2, 0, .5, .3)$

Verificăm dacă tip este una din literele “b”, “p” sau “n”. Dacă $tip = “b”$, avem de a face cu un contor de tip Binom(N, pN). Nu verificăm dacă N este număr natural, e răspunderea utilizatorului. Dacă $tip = “p”$ atunci contorul este repartizat Poisson(N). Nu mai verificăm pN . Dacă, în fine, $tip = “n”$, atunci $N \sim Negbin(N, pN)$. Este o repartiție negativă binomială generalizată, în sensul că N nu trebuie să fie neapărat număr natural.

Următorul pas este să calculăm ES și $Var(S)$

Aplicînd formulele lui Wald¹³ găsim

$$ES = EX \cdot EN, VarS = EN \cdot Var(X) + E^2N \cdot Var(N) \quad (6.3.9)$$

Pe baza lor calculăm dimensiunea lui f și F , anume $n_{\max} = ES + 15\sigma(S)$, $sue = 1 - \epsilon$

Acestea au fost pregătirile. Partea principală a scriptului este calculul recurenței (6.3.7). Cum recurența (6.3.7) este liniară în raport cu f_0 , preferăm să punem $f[1] = 1$ și să înmulțim la sfîrșit cu

$$f_0 = g(p_0)$$

¹³ Vezi de exemplu G. Zbăganu, op. Cit, p. 164

Scriptul este următorul:

```
##functia Panjer (x,p,"type",N,pN,epsilon) calculeaza repartitia lui S(N),
## N este un contor Panjer(a,b)
## daca type = "p", N ~ Poisson (N);se ia in considerare doar N, pN=0
## daca type = "b" e vorba de un contor binomial de parametru (N,pN)
## daca type = "n" e contor negativ binomial Bin(N,pN)
## 1. verificam daca x,p sunt de aceeasi lungime si daca p e probabilitate
## 2. Daca nu, mesaj de eroare, daca da, calculam a,b, f0
## 3. Calculam nmax = length(f) in asa fel incit sa evitam dimensiuni prea
mari
## nmax= int(ES + 10 sig(S))+1 cu ES = EXEN, sig^2(S)= ENVarX+VarNE^2X
## 4. Completam x ca sa fie de forma 0,1,...,(nmax-1) si p ca length(p)=nmax
## calculam f si F pina cind F*f0>(1-epsilon);
## Ca sa evitam zerourile de masina incepem cu f[1]=1
## recurenta este f[n+1]=sum(f=p[j+1]f[n+1-j] (a+bj/n), 1<=j<=n

Panjer<-function(x,p,type,N,pN,epsilon)
{
  x = floor(abs(x)); pN=min(pN,1) ## x are componente nr naturale si pN<1
  nx=length(x);np=length(p);
  if (nx!=np) {print ("Eroare: x si p au lungimi diferite");return()}
  for (i in 1:nx)
    {if (p[i]<0) {print("Eroare:p nu e probabilitate, are componente
negative");return()}}
    if (abs(sum(p)-1)>10^(-9)) {print("Eroare:p nu e probabilitate, suma nu e
1");return()}}
  ## standardizam pe x si p
  o=order(x); x=x[o];p=p[o] ## am aranjat ca x sa fie crescator
  n1=max(x)+1; x1=1:n1;x1=x1-1;p1=c(rep(0,n1))
  for (j in 1:nx) {p1[x[j]+1]=p[j]};sue=1-epsilon

  #3 initializam a,b,f0, ES, VarS
  mux=x%*%p;varx=(x^2)%*%p - mux^2;qN=1-pN
  if (type=="b")
  {muN = N*pN;varN=N*pN*qN; mus=mux*muN;vars=muN*varx+varN*mux^2;
  f0=(qN+pN*p1[1])^N;a=-pN/qN;b=-a*(N+1)}
  else if (type == "p")
  {muN=N;varN=N;mus=mux*muN;vars=muN*varx+varN*mux^2;f0=exp (N*(p1[1]-
1));a=0;b=N}
  else if (type == "n")
  {muN=N*qN/pN;varN=N*qN/pN/pN;mus=mux*muN;vars=muN*varx+varN*mux^2
  f0=(pN/(1-qN*p1[1]))^N;a=qN;b=a*(N-1)}
  else {print ("Nu inteleg tipul")}
  }
  #4 calculam nmax si initializam f,F
  nmax=mux+15*sqrt(vars);nmax=floor(nmax)+1;p=c(p1,rep(0,(nmax-n1)))
  f=c(rep(0,n1));F=f;x=x1;su=0
  ct=1/(1-a*p[1]);sue=(1-epsilon)/f0;n=1;f[1]=1;F[1]=1
  for (n in 1:nmax)
    {f[n+1]=0;n2=min(n,n1)
    for (j in 1:n2)
      { f[n+1]=f[n+1]+p[j+1]*f[n+1-j]*(a + b*j/n)}
    f[n+1]=ct*f[n+1]
    F[n+1]=F[n] + f[n+1]
    n=n+1
    if (F[n]>sue) break
  }
  f=f0*f;F=f0*F
  a=list(ES=mux,sigmaS=sqrt(vars),n_epsilon=n,f=f,F=F)}
}
```

Apelarea lui se face prin comanda

a = Panjer(x,p,type,N,pN,epsilon)

(3.10)

Functia furnizează o listă a cu

- a[[1]] = ES

- $a[[2]] = \sigma(S)$
- $a[[3]] = n_\varepsilon$, unde n_ε este primul n cu proprietatea că $P(S \leq n_\varepsilon) > 1 - \varepsilon$
- $a[[4]] = f$, cu $f[n] = P(S = n-1)$
- $a[[5]] = F$ cu $F[n] = P(S \leq n-1)$

Revenim acum la *Exemplul 1*: $X_n \sim \text{Bin}(10, \frac{1}{2})$ și $N \sim \text{Poisson}(100)$, în care ne așteptăm să avem în jur de 100 de clienți și fiecare să ceară, în medie 5 lei. Dacă acceptăm un risc de 1% , ne interesează $a[[3]]$ cu $\varepsilon = .01$. Probabilitățile p_j se calculează cu ajutorul funcției dbinom:

```
x=0:10; p=c(rep(0,11)); for (i in 1:11) {p[i]=dbinom(i-
1,10,.5)}
[1] 0.0009765625 0.0097656250 0.0439453125 0.1171875000 0.2050781250
[6] 0.2460937500 0.2050781250 0.1171875000 0.0439453125 0.0097656250
[11] 0.0009765625
```

Comanda principală este

```
a=Panjer(x,p,"p",100,.1,.01)
> ES=a[[1]];ES
[1,] 500
> sigmaS = a[[2]]; sigmaS
[1,] 52.44044
> VaR_epsilon= a[[3]]; VaR_epsilon
[1] 627
```

Deci în 99% din cazuri, suma cerută este mai mică decât 627 lei. Cu inegalitatea Cebîșev îmbunătățită găsisem 871 lei. Dacă dorim să vizualizăm densitatea f și funcția de repartiție F putem proceda astfel:

```
x=1:nn;plot(x,f/max(f),type="l",col="blue",sub="Albastru:
densitatea lui S, rosu: functia de repartitie");
lines(F,col="red") .
```

Obținem graficul

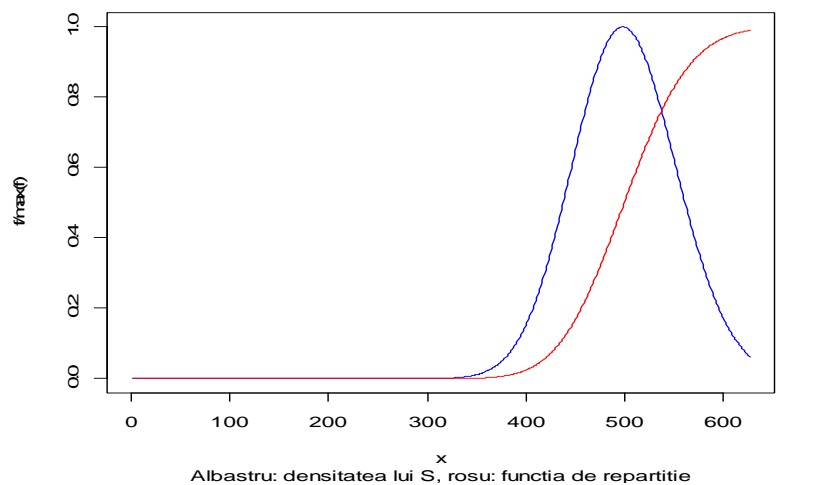


Fig 6.5

Se vede clar cum cele mai probabile valori sunt în jurul lui $x=500$. Dacă am fi vrut să ne asumăm riscul de 0.1% în loc de 1%, atunci valoarea la risc ar fi fost 671. La riscul de 0.01% valoarea ar fi fost $VaR = 708$ iar la $\epsilon = 10^{-6}$ ar fi fost 771 lei. Valoarea de 871 lei, găsită la riscul 1% cu inegalitatea lui Cebîșev ar fi suficientă la un risc de $\epsilon = 10^{-10}$!!

```
> a=Panjer(x,p,"p",100,1/11,.0000000001);VaR_S=a[[3]];VaR_S
[1] 872
```

Exemplul 3. Să zicem că $X_n \sim \text{Bin}(10, \frac{1}{2})$ și $N \sim \text{Negbin}(10, 1/11)$. Și acum ne așteptăm să avem în jur de 100 de clienți și fiecare să ceară, în medie 5 lei. Dar diferența este destul de mare

```
a=Panjer(x,p,"n",10,1/11,.01)
ES=a[[1]];sigS=a[[2]];VaR_S=a[[3]]; f = a[[4]]; F=a[[5]]
ES = 500
sigS = 166.5833
VaR_S = 964
```

Deși media este aceeași, acum abaterea medie pătratică este de trei ori mai mare ca înainte: 166.58. Valoarea la risc e normal să fie mai mare: acum este 964 lei. Graficul arată diferența

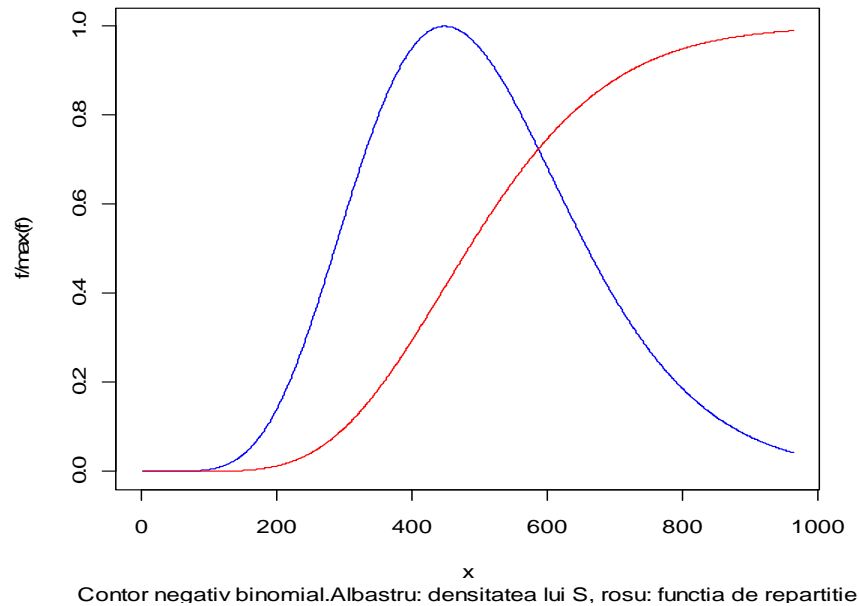


Fig 6.6

La riscul 0.1%, suma necesară ar fi fost 1168.

Putem acum compara algoritmul de convoluție $\text{convarep}(n,x,p)$ (algoritmul de bun simț) care calculează convoluția de n ori a repartiției $(x,p) = \begin{pmatrix} x_1 & x_2 & \dots & x_k \\ p_1 & p_2 & \dots & p_k \end{pmatrix}$ cu

algoritmul rapid convarap(x,p,N,ε) care calculează repartiția sumei $S = Y_1 + \dots + Y_N$ cu $N \sim \text{Bin}(n, 1 - P(X_n = 0))$, unde $Y_i = (X_i \mid X_i > 0)$ pînă valoarea la risc $\text{VaR}_{1-\varepsilon}(S)$ (adică pînă la primul n care are proprietatea că $P(S > n) < \varepsilon$.) Scriptul este următorul

```
Convarap <- function(x,p,n,epsilon)
{nx=length(x)
o=order(x);x=x[o];p=p[o] ##sortăm pe x
xmic=min(x);x1=x-xmic
## xmic devine 0, xmare devine xmare+xmic
pN=1-p[1];N=n
q=p/pN;q=q[2:nx]
y=x1[2:nx]; ## se rețin doar acei y>0
xa=Panjer(y,q,"b",N,pN,epsilon)
f=xa[[4]];F=xa[[5]];nn=length(f);x=0:(nn-1);x=x+N*xmic
b=list(x=x,f=f,F=F)
b
}
```

O primă observație este că algoritmul convarep permite efectuarea convoluțiilor fără nici o restricție, pe cînd la convarap trebuie ca $x_i \geq 0$.

O a doua este că n nu poate fi prea mare și nici $P(X = x_1)$ nu poate fi prea mic din următorul motiv: la inițializare, $f_0 = P(X = x_1)^n$.

Să luăm cazul unui zar. $x = (1,2,3,4,5,6)$, $p = (1,1,1,1,1,1)/6$. $x_1 = 1/6$, deci $P(S = n) = 1/6^n$. Numerele în „R” sunt date în dublă precizie, adică sunt cuprinse între 10^{-63} și 10^{63} . Dacă o să facem $n = 1000$, deja nu mai avem încredere în f_0 , care devine 0. Șirul (f_n) capătă un comportament aleator. Totuși, pentru $n = 400$, ambele scripturi sunt funcționale, al doilea fiind de 120 de ori mai rapid ca primul.

Pentru valori mici ale lui n și mai mari ale lui k , în mod paradoxal este mai bun algoritmul „baby”. De ex. dacă $n = 100$, $x = (1,3,5,7,11,13,17,19,23,29)$, $p = (1,1,1,1,1,1,1,1,1,1)/10$, atunci suntem mai siguri pe algoritmul naiv convarep, deși durează 30 de secunde, decît pe cel sofisticat, care durează jumătate de secundă

```
b=convarap(x,p,100,.00000001);t1=b[[1]];f1=b[[2]]
a=convarep(100,x,p);t=a[[1]];f=a[[2]]
plot(t1,f1,type="l");lines(t,f,col="red")
```

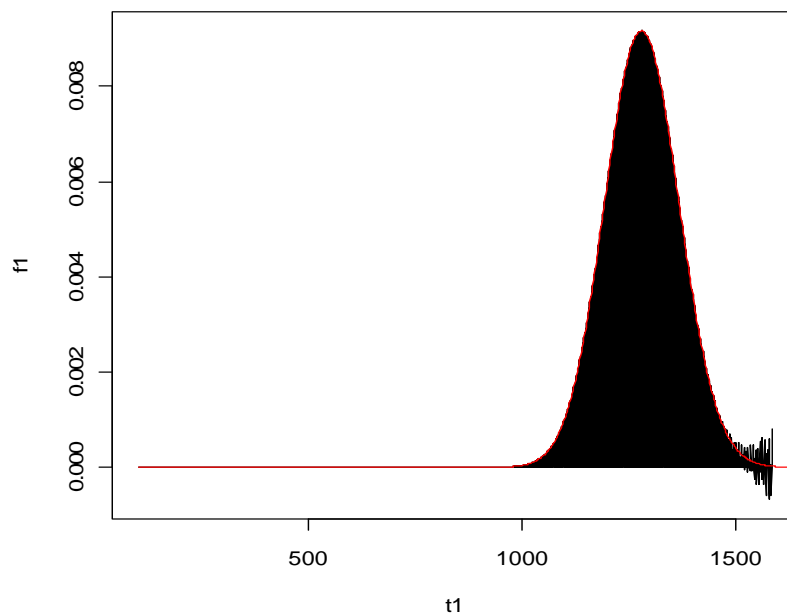



Fig 6.7

Cu negru sunt figurate densitățile f_i calculate după algoritmul rapid. Deja unele devin negative.

Există alte tehnici sofisticate de a rezolva problema convoluției dar nu fac obiectul lucrării de față.

6.4 . Probleme de demografie, asigurări de persoane, tabele de mortalitate

Un tabel de mortalitate este un tabel cu două coloane, x și l în care x reprezintă vârsta iar l_x este probabilitatea ca un individ născut astăzi să depășească vârsta de x ani, înmulțită cu 100.000. Deci neapărat $l_0 = 100000$. Tabelele de mortalitate se construiesc pînă la o vîrstă maximă, notată convențional cu ω . De obicei $\omega = 100$ sau $\omega = 110$.

Datele brute pe baza cărora se construiește tabelul sunt doi vectori, \mathbf{n} și \mathbf{d} unde n_i este numărul total al populației cu vîrstă cuprinsă în intervalul $[i, i+1)$ și d_i numărul de decese ale populației cu vîrste cuprinse între aceste limite.

De exemplu, n_0 este numărul de copii cu vîrste cuprinse între 0 și 1 an, calculat la 1 ianuarie iar d_0 este numărul deceselor acestei categorii de populație, calculat în același an, dar la 31 decembrie. Sunt și alte convenții, dar aceasta este cea mai simplă.

Iată un **exemplu** de date brute, luat de la institutul de statistică în care avem de a face cu patru populații: bărbați 1994, femei 1994, bărbați 2008 și femei 2008. Exemplul este pus într-un fișier .txt cu numele de `barbatfemei19942008` care trebuie să fie în directorul de lucru, pentru a nu fi complicat de deschis cu instrucțiunea `read.table`.

x	1994				2008			
	barbati total	decese	femei total	decese	barbati total	decese	femei total	
decese								
0	127308	3303	120921	2591	110665	1416	105025	1018
1	125979	282	119583	249	109056	93	103456	77
2	131763	197	125846	132	112199	65	105912	49
3	127804	158	122161	110	110315	52	104550	37
4	175284	192	167954	124	107839	34	101622	24
5	176293	255	170667	159	106203	40	100940	37
6	182212	179	173773	140	107077	42	100304	30
7	188160	118	181410	80	109577	46	104001	28
8	169580	95	163040	54	114464	35	108485	20
9	175287	103	167986	43	113859	30	108823	17
10	157988	91	151287	44	114695	31	108680	13
...								
80	29268	3576	48002	4847	40991	4087	65403	4968
81	28345	3743	46395	4975	35451	3848	58600	4947
82	24626	3490	40820	4789	31021	3585	52227	4992
83	20438	3179	33543	4396	26196	3498	45978	4981
84	17178	2930	28381	4250	21888	3082	39356	4833
85>	57985	13476	98009	21478	66841	13225	133181	25154

De exemplu, dacă ne uităm la $x = 80$, vedem că la 1 ianuarie 1994 erau 29268 bărbați și 48002 femei cu vârste între 80-81 de ani, iar la 31 decembrie 1994 muriseră din ei 3576 bărbați și 4847 femei. Tot în acel rând vedem că la 1 ianuarie 2008 erau 40991 bărbați și 65403 femei din care, la 31 decembrie 2008 decedaseră 4087 bărbați și 4968 femei.

Se poate vedea că ultimul rând face excepție: aici sunt cuprinse toate persoanele cu vârste de peste 85 de ani. La 1 ianuarie 1994 erau 57985 și 98009 femei (din ei decedaseră la 31 decembrie 2008 un număr de 13476 bărbați și 21478 femei) iar la 1 ianuarie 2008 erau 66841 bărbați și 133181 femei din care la 31 decembrie decedaseră 13225 bărbați și 25154 femei.

Completarea tabelului pînă la vârsta de 85 de ani se face identificînd probabilitatea unui eveniment cu frecvența lui.

Mai precis: să presupunem că timpul de viață la naștere al unui individ este o variabilă aleatoare T . Vrem să calculăm coada sa, $\underline{F}(x) = P(T > x)$, care în demografie se notează cu ${}_x p_0$. Normal ar trebui să punem ${}_x p_0 = \frac{\text{total supraviețuitori cu vârste} > x}{\text{total supraviețuitori}}$,

dar acest lucru este imposibil, deoarece nu știm numărul supraviețuitorilor: ar trebui să știm vârsta decesului...

Dacă acceptăm identificarea probabilității cu frecvența, interpretăm raportul $p_0 = d_0/n_0$ cu probabilitatea $P(T \leq 1)$: probabilitatea ca un nou născut să moară înainte de a împlini un an. Atunci numărul $q_0 = 1 - p_0$ este $P(T > 1)$, adică probabilitatea ca el să supraviețuiască un an. Apoi, $p_1 = d_1/n_1$ este probabilitatea ca un supraviețuitor de 1 an să moară în intervalul de timp $(1, 2]$, adică înseamnă $P(T \leq 2 \mid T > 1)$. Deci $q_1 = 1 - p_1$ înseamnă probabilitatea ca un individ să supraviețuiască 2 ani știind că el deja a supraviețuit 1 an, adică $q_1 = P(T > 2 \mid T > 1)$, adică avem $q_1 = P(T > 2)/P(T > 1) \Rightarrow P(T > 2) = q_0 q_1$. Din aproape în aproape găsim formula

$$P(T > x+1) = q_1 \dots q_x \quad (6.4.1)$$

valabilă pentru x număr natural. Înmulțim acest număr cu 100000 și obținem $l_x = 10^5 q_x$, rotunjit la cel mai apropiat întreg. Rezultatul este valabil, în cazul exemplului nostru numai pînă la $x = 84$, deoarece nu avem date mai multe.

Funcția `tabmort` realizează acest lucru.

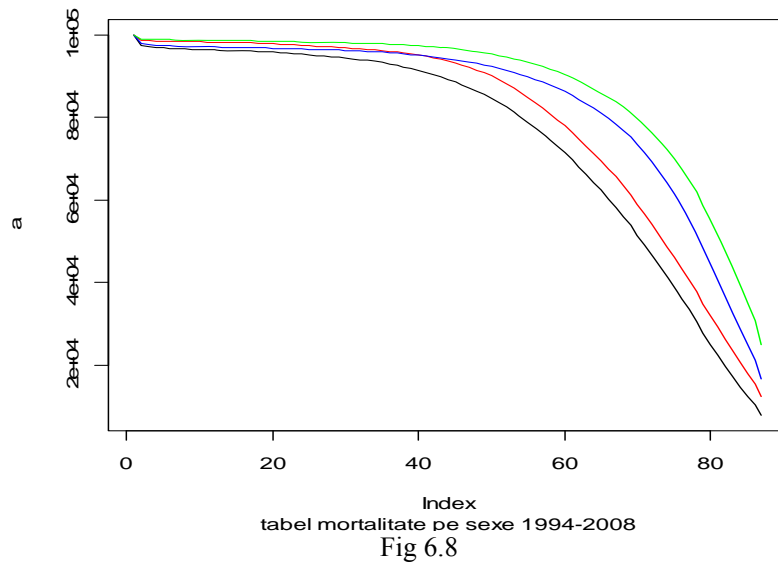
```
## funcția tabmort (numar,death) creează un vector
## cu același număr de componente de elemente
tabmort<-function(numar,death)
{nv=length(numar);nm=length(death);tabel=c(rep(0,nv))
if (nv!=nm) {print ("eroare:cei doi vectori au lungimi
diferite");return() }
p=death/numar;q=1-p
for (j in 1:nv) {tabel[j]=prod(q[1:j])}
tabel=c(1,tabel)
tabel = floor(100000*tabel+.5)
tabel
}
```

Pentru a o aplica, am făcut o pregătire: am copiat în clipboard fișierul `barbatfemei19942008` (cu excepția primelor două rînduri, care nu se potrivesc cu restul) și am obținut o listă cu 9 elemente, numită `tabgen`

```
tabgen=read.table("clipboard")
```

Din lista `tabgen` extragem cei 8 vectori (`numar,death`) cu care creăm apoi 4 tabele de mortalitate, notate cu `a,b,c,d` pe care le plotăm. Primul se referă la bărbați 1994 (linia neagră) al doilea la bărbați 2008 (cea roșie), al treilea la femei 1994 (albastră) și ultimul la femei 2008 (verde)

```
> virsta=tabgen[[1]];nrbb94=tabgen[[2]]
nrbb94=tabgen[[3]];nrfm94=tabgen[[4]];nrmm94=tabgen[[5]]
> nrbb08=tabgen[[6]];nrbb08=tabgen[[7]]
nrfm08=tabgen[[8]];nrmm08=tabgen[[9]]
a=tabmort(nrbb94,nrbb94);plot(a,type="l",sub="tabel
mortalitate pe sexe 1994-2008" )
b=tabmort(nrbb08,nrbb08);lines(b,col="red")
c=tabmort(nrfm94,nrmm94);lines(c,col="blue")
d=tabmort(nrfm08,nrmm08);lines(d,col="green")
```



Pentru a putea completa tabelul pînă la $\omega = 100$ sau $\omega = 110$ trebuie să acceptăm unele ipoteze neverificabile. Una din ele ar fi că **riscul de moarte** este constant pe intervale de forma $(k, k+1]$ și pe intervalul $(84, \infty)$ crește liniar. Alta ar fi că riscul de moarte crește exponențial (ipoteza lui Gompertz).

În ambele modele se ascunde o ipoteză neverificabilă, anume că repartitia timpului mediu de viață la naștere, T este absolut continuă și coada sa, $\underline{F}_T(x)$ este de forma

$$\underline{F}_T(x) = e^{-\int_0^x r(y)dy} \Leftrightarrow r(x) = -\frac{F_T'(x)}{F_T(x)} \quad (6.4.2)$$

Funcția r se numește **riscul instantaneu de moarte**. Dacă acceptăm că r este o funcție în scară, constantă pe intervale de forma $(k, k+1)$, atunci avem că ${}_{x+1}p_0 = {}_xp_0 e^{-\int_x^{x+1} r(y)dy}$ de unde deducem expresia

$$r_x = \ln \frac{l_x}{l_{x+1}} \Leftrightarrow e^{r_x} = \frac{l_x}{l_{x+1}} \quad (6.4.3)$$

În figura de mai jos prezentăm riscurile de moarte la cele patru populații: bărbați 1994/2008 (curba neagră și cea roșie) și femei 1994/2008 (curba albastră și verde)

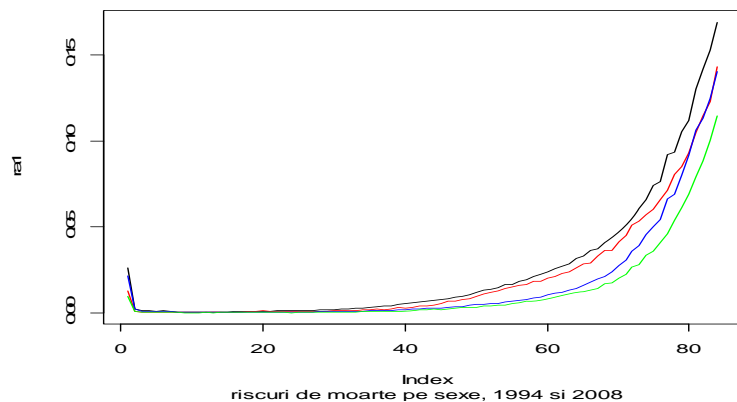


Fig. 6.9

Figura sugerează că ar fi acceptabil un model în care să se presupună că de la 85 de ani încolo, riscul de moarte r crește liniar. Pe această ipoteză este construit scriptul $tp0(t, tab)$ care calculează funcția de supraviețuire ${}_tp_0$ pentru orice t plecând de la un tabel de mortalitate, tab . Din tabelul de mortalitate l reținem doar acele poziții x în care $l_x > 10$. Atunci, din formula (4.3) deducem că

$${}_tp_0 = l_{[t]} \left(\frac{l_{[t]+1}}{l_{[t]}} \right)^{\{t\}} \text{ dacă } t \leq \omega \text{ și } {}_tp_0 = l_{\omega} q^{\frac{(t-\omega)(t-\omega+1)}{2}} \left(\frac{q}{p} \right) \text{ dacă } t > \omega \quad (6.4.4)$$

Am notat $p = \frac{l_{\omega-1}}{l_{\omega-2}}$ și $q = \frac{l_{\omega}}{l_{\omega-1}}$.

```
## Functia tp0 (t, tabel de mortalitate) calculeaza functia
## de supravietuire a unui nou nascut la momentul t in ipoteza ca
## riscul de moarte creste liniar dupa virsta nv. Ca sa evitam
## impartiri cu 0, retinem din tabelul de mortalitate "tab" doar
## elementele mai mari ca 10
tp0<-function(t,tab)
{nv=length(tab);ntv=n-1;k=floor(t);s=t-k
nvv=length(which(tab>10))
p=tab[nvv-1]/tab[nvv-2];q=tab[nvv]/tab[nvv-1]
r=q/p;ss=t-nvv+1
if (k<(nvv-1)) {tp0=tab[k+1]*(tab[k+2]/tab[k+1])^s}
else tp0=tab[nvv]*q^(ss)*r^((ss^2+ss)/2)
tp0=tp0/100000
#list(tp0=tp0,nv=nv,r=r,k=k)
tp0}
```

Cu ajutorul funcției care calculează ${}_tp_0$ pentru orice număr pozitiv, nu neapărat întreg, este ușor de calculat funcția de supraviețuire la vârsta x , care se notează cu ${}_xp_x$ și se definește ca

$${}_xp_x = P(T-x > t | T > x) = \frac{P(T > x+t)}{P(T > x)} = \frac{{}_{x+t}p_0}{{}_xp_0} \quad (6.4.5)$$

Scriptul este

```
## functia tpx (t,x,tabel de mortalitate) calculeaza probabilitatea
## de supravietuire a unui individ de x ani in aceeași ipoteza ca și
## tp0
tpx<-function(x,t,tab)
{tpx=tp0(t+x,tab)/tp0(x,tab);tpx}
```

Variabila aleatoare cu repartiția (4.5) se notează cu T_x : timpul rezidual de viață la vârsta de x ani. Deci $T_x = (T - x | T > x)$. Media acestei variabile aleatoare este notată cu e_x . Aplicînd formula $ET_x = \int_0^{\infty} P(T_x > t) dt = \int_0^{\infty} {}_tp_x dt$, pe care o aproximăm cu $e_x = \sum_{x \geq 0} {}_tp_x$ formăm scriptul

```
## speranta medie de viata a unui individ in virsta de x ani
ex<-function(x,tab){ex=0
for (j in 0:120) {ex=ex+tpx(x,j,tab)}
ex}
```

În figura de mai jos este prezentat graficul funcției ex , respectând culorile din celelalte două figure de mai sus: negru înseamnă “bărbați 1994”, roșu înseamnă “bărbați 2008”, albastru este “femei 1994” iar verde este “femei 2008”

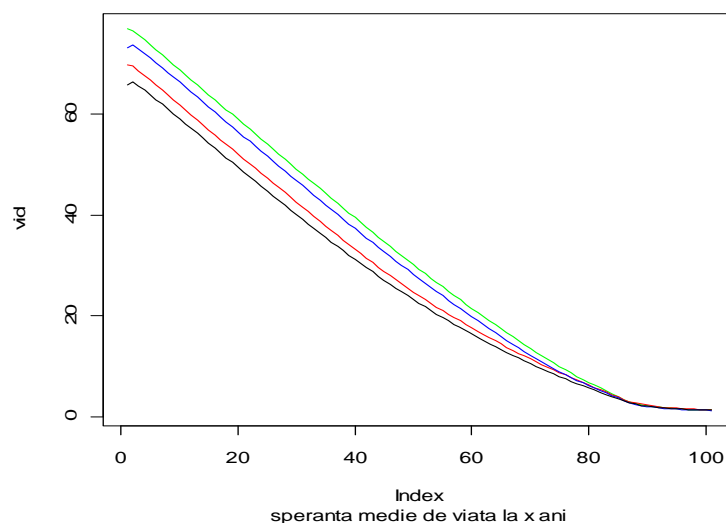


Fig 6.10

Tabelele acestea sunt importante în studiul asigurărilor de viață.

Iată un exemplu tipic.

Un individ de vîrstă x vrea să își facă următoarea asigurare de viață: să primească o sumă de S_0 lei peste t ani dacă va fi în viață, iar dacă nu, urmașii lui să primească aceeași sumă la momentul decesului. Pentru aceasta el constituie un depozit bancar cu rata anuală a dobînzii i la care să aibă dreptul el (numai peste n ani) sau moștenitorii lui (la momentul T_x al decesului)

Ce sumă Π să depună în așa fel încît probabilitatea ca suma să fie insuficientă să fie ε ?

Datele problemei sunt $S_0, x, n, i, \varepsilon$. Formal, abordarea este următoarea: fie T_x timpul rezidual de viață al individului și $\tau = \min(T_x, n)$. Dacă suma depusă este Π , atunci suma pe care o va avea el sau urmașii lui peste τ ani va fi $S_1 = \Pi(1+i)^\tau$. Condiția pusă este ca $S_1 > S_0$ și se cere cel mai mic Π cu proprietatea că $P(S_1 > S_0) \geq 1 - \varepsilon$.

$$\text{Cum } S_1 > S_0 \Leftrightarrow (1+i)^\tau > \frac{S_0}{\Pi} \Leftrightarrow \tau \ln(1+i) > \ln \frac{S_0}{\Pi} \Leftrightarrow \tau > \frac{\ln \frac{S_0}{\Pi}}{\ln(1+i)} \text{ rezultă că}$$

suma minimă care trebuie depusă, Π va avea valoarea

$$\Pi = S_0(1+i)^{-y} \text{ unde } y \text{ este soluția ecuației } P(\tau > y) = 1 - \varepsilon \quad (6.4.6)$$

În concluzie, problema se reduce la calculul cuantilelor repartiției lui τ . Cuantila este inversa funcției de repartiție. Altfel spus, este soluția y a ecuației se numește cuantila la nivel ε a lui τ . În cazul nostru avem

$$P(\tau > y) = 1 - \varepsilon \Leftrightarrow P(\min(T_x, n) > y) = 1 - \varepsilon \Leftrightarrow P(T_x > y) = 1 - \varepsilon, y < n \Leftrightarrow y p_x = 1 - \varepsilon$$

Deci y este cuantila $F_\tau^{-1}(\varepsilon)$

```
##cuantilele lui tpx
cuant<-function(x,tab,p){p=1-
p;pas=.05;sa=seq(0,120,by=pas);pa=sa;na=length(sa)
for (j in 1:na) {pa[j]=tpx(x,sa[j],tab)}
v=which(pa>p);ja=length(v);ja1=ja+1;alfa=pa[ja];beta=pa[ja
1];tot=alfa+beta
y=beta*ja/tot+alfa*ja1/tot;y=y*pas
y}
```

Atunci suma Π care trebuie depusă în depozit pentru a garanta cu riscul ε existența unei sume S_0 la momentul $\tau = \min(T_x, n)$ va fi

$$\Pi = S_0 / (1 + i)^{\text{cuant}(x, \text{tab}, p)}$$

De exemplu, pentru a afla suma necesară unui bărbat de 60 de ani care dorește ca la momentul τ să aibă în depozit suma de 10.000 RON dacă dobânda este de 5% și folosim tabelul de mortalitate pentru 2008, scriem comenzile

```
b=tabmort(nrbb08,nrbm08);
y=cuant(60,b,.05);PI=1000*1.05^(-y);PI
[1] 8927.597
```

Sau, mai profesionist

```
s0=10000;i=.05;PI=s0/(1+i)^y;PI
[1] 8927.597
```

Sau, ca să arătăm cum se poate face totul cu o singură comandă

```
PI=s0/(1+i)^cuant(60,tabmort(nrbb08,nrbm08),.05);PI
[1] 8927.597
```

Dacă dobânda era 7% și riscul de 1%, am fi avut

```
PI=s0/1.1^cuant(60,tabmort(nrbb08,nrbm08),.01);PI
[1] 9511.922
```

Dacă o femeie de 60 de ani vrea să facă același lucru, atunci am fi pus

```
PI=s0/1.1^cuant(60,tabmort(nrffm08,nrffm08),.01);PI
[1] 8983.24
```

Se poate vedea cum contează diferența de vîrstă și sex. Dacă modificăm vîrsta la 50 de ani, avem

```
PI=s0/1.1^cuant(50,tabmort(nrffm08,nrffm08),.01);PI = 7936.383
(la femei)
```

```
PI=s0/1.1^cuant(50,tabmort(nrbb08,nrbm08),.01);PI = 9112.59
(la bărbați)
```

Dacă vîrsta este de 5 ani, iar $S_0 = 40.000$ (asigurare de tip *endowing* – adică înzestrare) avem

$PI = s_0 / 1.1^{\text{cuant}(5, \text{tabmort}(\text{nr fm08}, \text{nr mm08}), .01)}$; $PI = 2600.906$
(la femei)

$PI = s_0 / 1.1^{\text{cuant}(5, \text{tabmort}(\text{nr bb08}, \text{nr bm08}), .01)}$; $PI = 7456.137$
(la bărbați)

Dacă doritorul de asigurare vrea să micșoreze costul asigurării și să păstreze același risc ca suma să fie suficientă, atunci ar trebui să se asocieze cu alți asigurați de același tip, sau să apeleze la un asigurator.

Ideea (pe care o vom explica pe primul exemplu, al bărbatului de 60 de ani) este următoare : dacă actualizăm suma pe care o depune el ca la momentul $\min(n, T_x)$ să se obțină suma $S_0 = 10.000$ RON este o variabilă aleatoare, anume $X = S_0 / (1+i)^{\min(T_x, n)}$. Ar vrea să depună o sumă, Π cât mai mică cu putință ca $P(\Pi(1+i)^{\min(T_x, n)} < S_0) \leq \varepsilon$, adică $P(X > \Pi) \leq \varepsilon$. Dacă este singur, atunci va trebui să aleagă pe Π ca fiind cuantila de nivel $1 - \varepsilon$ a variabilei aleatoare X .

Dar dacă se asociază cu alții, în aceeași situație și depun cu toții aceeași sumă, Π ? Să presupunem că s-a format o asociație din N asemenea indivizi, dornici să coopereze. Să acceptăm că timpii lor reziduali, notați cu T_j au toți aceeași repartiție ca T_x . Să mai presupunem că T_j sunt independenți și că toți indivizii noștri fac un depozit comun în care depun **aceeași sumă, Π** . Fiecare își retrage bani de acolo atunci când este nevoie. Acesta este numit în asigurări **principiul solidarității**.

Fie $X_j = S_0 / (1+i)^{\min(T_j, n)}$. Variabilele aleatoare X_j vor fi de asemenea independente. Probabilitatea ca depozitul să fie insuficient este $P(X_1 + \dots + X_N > N\Pi)$. Dacă dorim ca această probabilitate să fie mai mică sau egală cu ε , va trebui să alegem Π în așa fel încât $n\Pi$ să fie cuantila la nivel $1 - \varepsilon$ a repartiției sumei $S_N = X_1 + \dots + X_N$, adică a convoluției F^{*N} , unde F este repartiția lui X_j .

Este complicat de lucrat cu convoluțiile, dacă repartiția este continuă. Dar se poate aplica inegalitatea lui Cebîșev. Să spunem că $EX_j = \mu$, $\sigma(X_j) = \sigma$. Aproximăm $E(X - \mu)_+^2$ cu $\sigma^2/2$. Atunci putem aplica inegalitatea lui Cebîșev de la (3.5) sub forma $P(S_N > ES_N + k\sigma(S_N)/\sqrt{2}) < 1/k^2$, care în cazul nostru devine $P(S_N > N\mu + k\sigma\sqrt{N/2}) < 1/k^2$. Pentru un risc de $\varepsilon = 4\%$, de exemplu, alegem $k = 5$ și avem $P(S_N > N\mu + 5\sigma\sqrt{N/2}) < .04$. Luăm Π în așa fel încât $N\Pi = \mu + 5\sigma\sqrt{N/2}$, adică $\Pi = \mu + \frac{5\sigma}{\sqrt{2N}}$. Pentru aceasta avem de calculat EX și $\sigma(X)$.

Dacă scriem $1+i = e^\delta$, atunci $X_j = S_0 e^{-\delta \min(T_j, n)}$, de unde $EX_j = S_0 L_n(\delta)$, $EX_j^2 = S_0^2 L_n(2\delta)$, unde

$$L_n(\delta) = E e^{-\delta \min(T_j, n)} \quad (6.4.7)$$

este transformata Laplace a variabilei aleatoare $\min(T_j, n)$. Aceasta se poate calcula foarte comod folosind formula de integrare prin părți

$$Eg(X) = g(0) + \int_0^\infty g'(t) \underline{F}_X(t) dt \quad (6.4.8)$$

pentru funcția $g(x) = e^{-\delta \min(x, n)}$. Cum $g(0) = 1$ și $g'(x) = -\delta e^{-\delta x} 1_{(0, n)}(x)$, formula de mai sus devine

$$L_n(\delta) = 1 - \delta \int_0^n e^{-\delta t} \underline{F}(t) dt, \text{ unde } \underline{F} \text{ este funcția de supraviețuire a variabilei aleatoare } T_x$$

(coada sa), care se notează cu p_x și se calculează cu ajutorul funcției

$tpx(x, t, tab)$ cu care am făcut deja cunoștință. În definitiv din (4.8) deducem formulele

$$L_n(\delta) = 1 - \delta \int_0^n e^{-\delta t} \cdot {}_t p_x dt, EX_j = S_0 L_n(\delta), \sigma(X_j) = S_0 \sqrt{L_n(2\delta) - L_n^2 \delta} \quad (6.4.9)$$

Scriptul `Laplace(x, n, delta, tab, nrint)` calculează pe $L_n(\delta)$ cu $\delta = \ln(1+i)$ prin metoda trapezelor, cu pasul diviziunii intervalului $[0, n]$ egal cu $n/nrint$.

Amintim semnificația parametrilor:

- x este vârsta asiguratului;
- $\delta = \ln(1+i)$ unde i este rata anuală a dobânzii;
- n este durata depozitului, dacă asiguratul nu moare în această perioadă;
- tab este tabelul de mortalitate pe baza căruia se face calculul
- $nrint$ este numărul intervalelor echidistante în care se împarte intervalul $[0, n]$ pentru a putea calcula integrala

```
## Laplace(x,n,delta,tab,nrint) calculeaza integrala 1 -
#delta int(exp(-tdelta) tpx
Laplace<-function(x,n,delta,tab,nrint)
{ pas=n/nrint;s=1;
for (j in 1:nrint) {s=s+exp(-delta*j*pas)*tpx(x,j*pas,tab)}
s=s-(1+exp(-delta*n)*tpx(x,n,tab))/2; s=s*pas*delta
rez=1-s;rez}
```

iar scriptul calculează media și varianța variabilei aleatoare $X = X_j = S_0 e^{-\delta \min(T_j, n)}$

```
##media si varianta daca x= 60,n=20, i=.05,
medvarxni<-function(s0,x,n,i,tab)
{
delta=log(1+i);nrint=12*n;la1=Laplace(x,n,delta,tab,nrint)
la2=Laplace(x,n,2*delta,tab,nrint)
ex=s0*la1;sigx=sqrt(la2-la1^2)
rez=list(EX=ex,SigmaX=s0*sigx)
rez }
```

Revenim la individul de 60 de ani care dorește 10000 RON peste 20 de ani. Calculul pe baza tablei din 1994 ne dă

```
> medvarxni(s0,x,n,i,a)
$EX 5374.11
$SigmaX 1796.007
```

Pe baza tablei din 2008 avem

```
medvarxni(s0,x,n,i,b)
$EX] 5188.782
$SigmaX 1738.958
```

Dacă e vorba de o femeie, atunci tabela din 1994 dă rezultatul

```
> medvarxni(s0,x,n,i,c)
$EX 4754.943
$SigmaX 1481.948
```

iar dacă folosim tabela din 2008,

```
> medvarxni(s0,x,n,i,d)
$EX 4549.405
$SigmaX 1377.932
```

Marginea $\Pi = \mu + \frac{5\sigma}{\sqrt{2N}}$, care garantează un risc de 4% de insuficiență a depozitului devine $\Pi_N = 5374 + \frac{6350}{\sqrt{N}}$.

Comparată cu valoarea exactă, în cazul unei autoasigurări cu risc 4% care este $\Pi_0 = 9238$, ea devine deja mai mică dacă numărul participanților la asigurarea în grup este mai mare ca 3. Într-adevăr, $\Pi_2 = 9864$, dar $\Pi_3 = 9040 < 9238$, $\Pi_4 = 8549, \dots, \Pi_{100} = 6009$.

Aceleași rezultate se obțin și dacă folosim tabele din 2008. Autoasigurarea realizată de

```
PI=s0/(1+i)^cuant(60,tabmort(nrb08,nrbm08),.04);PI
```

ne dă $\Pi_0 = 9104$, iar formula $\Pi_N = \mu + \frac{5\sigma}{\sqrt{2N}}$ ne dă $\Pi_3 = 8738 < 9104$, $\Pi_4 = 8263, \dots, \Pi_{100} = 5804$

Spre deosebire de cazul autoasigurării, fie în individual, fie în grup, în cazul în care se apelează la o companie de asigurări (care din asemenea lucruri trăiește), la primele cu risc de forma $EX + \alpha\sigma(X)$ se mai adaugă și un comision. Chiar și așa, este o soluție preferabilă față de autoasigurarea individuală.

6.5. Simulări de variabile aleatoare, probleme de portofolii

După cum am văzut la capitolul anterior, limbajul „R” oferă multe repartiții cărora li s ecalculează funcția de repartiție (p) densitatea (d), cuantila (q) și se pot simula (r). Ele sunt

Repartiția	Numele ei în R	Parametri
beta	beta	α, β
binomială	binom	k, p
Cauchy	cauchy	m, a
χ^2	chisQ	n
exponentială	exp	λ
F	f	m, n
gamma	gamma	v, λ
geometrică	geom	p
hipergeometrică	hyper	a, n, k^{14}
log-normală	lnorm	μ, σ
logistică	logis	μ, σ^{15}

¹⁴ Extragem k bile dintr-o urnă cu a bile albe și n bile negre; X este numărul de bile albe.

¹⁵ Repartiția logistică are funcția de repartiție $\frac{1}{1 + e^{-\frac{x-\mu}{\sigma}}}$. Este mai rar folosită..

negativ binomială	nbinom	k, p
normală	norm	μ, σ^{16}
Poisson	pois	λ
Student	t	n
uniformă	unif	a, b
Weibull	weibull	k, λ
Wilcoxon	wilcox	m, n

Pentru a genera un vector cu n componente i.i.d. cu aceste repartiții se pune în fața numelui lor din “R” litera r . Apoi se pune numărul de variabile aleatoare dorite și parametrii repartiției.

De exemplu:

```
x=rnorm(100,10,4);x
# x este un vector cu 100 de componente repartizate  $N(10,4^2)$ 
x=rbinom(10,10,.4);x
[1] 3 5 1 5 4 6 4 5 6 3
# x este un vector cu 10 #componente repartizate Binomial(10,.4)
x=rnbinom(6,10,.4);x
[1] 19 8 9 25 8 19
# x este un vector cu 10 componente repartizate Negbin(10,.4)
x=rhyper(10,4,6,6);x
[1] 3 2 3 3 2 2 2 2 1 3
# x este un vector cu 10 componente ~Hypergeometric(4,6,6)
```

Problemă. Cum se simulează o repartiție care nu face parte dintre acestea?
De exemplu, cum simulăm timpii de viață reziduali de viață $(T_j)_{1 \leq j \leq n}$ ai unor indivizi care au aceeași vîrstă, x ?

Dacă este vorba de repartiții discrete, limbajul “R” pune la dispoziție două instrucțiuni, care fac, printre altele acest lucru .

```
sample(x, size, replace = FALSE/TRUE, prob = NULL/p)
sample.int(n, size = n, replace = FALSE, prob = NULL)
```

Aici x este un vector în care se pun valorile variabilei aleatoare X (nu neapărat numere, pot fi și simboluri); p este un vector de aceeași lungime ca și x format cu numere pozitive, din care se formează automat o probabilitate, de către “R”, înlocuind $p[j]$ cu $p[j]/\text{sum}(p)$. Dacă nu se pune nici o probabilitate, p , se consideră că p este repartiția uniformă, $p = c(\text{rep}(1, \text{length}(x)))$. Dacă parametrul “replace” (sau chiar “rep”, se accepta prescurtaqrea!) este TRUE, atunci instrucțiunea `sample(x, n, rep=TRUE, prob=p)` simulează un șir de n variabile aleatoare $X_j \sim \begin{pmatrix} x_1 & x_2 & \dots & x_k \\ p_1 & p_2 & \dots & p_k \end{pmatrix}$, cu $k = \text{length}(x)$.

Exemplu. Să se simuleze un șir de 100 variabile aleatoare $X_j \sim \begin{pmatrix} a & b & c & d \\ .1 & .2 & .3 & .4 \end{pmatrix}$.

Soluție.

```
x= c("a", "b", "c", "d"); nr=100; xsim=sample(x,nr,rep=TRUE,prob=1:4)
> xsim
```

¹⁶ În codificarea “R”, `x=rnorm(1,μ,σ)` produce o variabilă aleatoare $X \sim N(\mu, \sigma^2)$. Parametrul al doilea reprezintă abaterea medie pătratică și nu varianța. Uneori se face confuzie.

```

[1] "a" "b" "c" "d" "d" "a" "d" "d" "b" "a" "c" "a" "d" "b" "c" "d" "c"
"d"
[19] "c" "d" "c" "d" "c" "c" "c" "b" "b" "d" "c" "d" "b" "c" "d" "c" "b"
"a"
[37] "c" "d" "c" "a" "b" "b" "c" "a" "c" "b" "b" "a" "a" "b" "a" "c" "c"
"c"
[55] "d" "d" "b" "d" "b" "c" "d" "d" "c" "d" "a" "d" "a" "d" "a" "c" "d"
"a"
[73] "c" "b" "d" "d" "c" "c" "c" "c" "d" "c" "d" "a" "d" "a" "d" "b" "c"
"d"
[91] "c" "d" "c" "a" "c" "d" "c" "b" "d" "b"

```

Ne putem convinge că este așa folosind instrucțiunea `table`, care ne arată frecvențele empirice ale celor patru simboluri

```

> table(xsim)
xsim
 a  b  c  d
17 18 33 32

```

Dacă simulam mai multe, se vedea mai clar cum probabilitățile sunt proportionale cu 1,2,3,4. de exemplu, dacă scriem

```

nr=10000; xsim=sample(x,nr,rep=TRUE,prob=1:4); table(xsim)
xsim
 a    b    c    d
1011 1983 3000 4006

```

acest lucru se vede mai clar.

Instrucțiunea

`sample.int(n, size = n, rep = FALSE/TRUE, prob = NULL/p)`
face cam același lucru, în cazul particular în care $x = 1:k$. De exemplu, dacă în loc de “a”, “b”, “c”, “d” am fi pus 1,2,3,4 , am fi putut folosi, cu același rezultat instrucțiunea

```

xsim=sample.int(4,nr,rep=TRUE,prob=1:4); table(xsim)
xsim
 1    2    3    4
1020 1997 2870 4113

```

Dar instrucțiunea `sample` face mai mult de atât: dacă `rep = FALSE`, atunci generează un aranjament aleator al componentelor lui x . Dacă `prob=NULL`, adică nu punem nici o probabilitate, toate aranjamentele au aceeași șansă de apariție. Este ca și cum am avea o urnă cu bilele $x = (x_1, \dots, x_k)$ din care scoatem, succesiv și fără înlocuire, un număr de `size` bile. Dacă parametrul `size` lipsește, se generează o permutare aleatoare a lui x .

Exemple:

```

> x= c("a", "b", "c", "d");
nr=3;xsim=sample(x,nr,rep=FALSE); xsim
[1] "b" "c" "a"
xsim=sample(x,rep=FALSE); xsim
[1] "c" "b" "d" "a"
> xsim=sample.int(4,rep=FALSE); xsim
[1] 1 3 4 2

```

```
xsim=sample.int(4,2,rep=FALSE); xsim
[1] 3 1
```

Dacă există și parametrul p – adică probabilitatea, atunci

```
xsim=sample(x,size,rep=FALSE,prob=p)
```

generează un aranjament $x_{\sigma(1)}, \dots, x_{\sigma(\text{size})}$ în care probabilitățile p se aplică succesiv. Mai precis, dacă $X \sim \begin{pmatrix} x_1 & x_2 & \dots & x_k \\ p_1 & p_2 & \dots & p_k \end{pmatrix}$ și $\text{size} = d \leq k$, se generează un vector aleator $V = (V_1, \dots, V_d)$ cu repartiția

$$P(V_1 = x_{\sigma(1)}, V_2 = x_{\sigma(2)}, \dots, V_d = x_{\sigma(d)}) = p_{\sigma(1)} \frac{p_{\sigma(2)}}{1 - p_{\sigma(1)}} \frac{p_{\sigma(3)}}{1 - p_{\sigma(1)} - p_{\sigma(2)}} \dots \frac{p_{\sigma(d)}}{1 - p_{\sigma(1)} - \dots - p_{\sigma(d-1)}}$$

Exemplu. $X \sim \begin{pmatrix} 1 & 2 & 3 & 4 \\ p_1 & p_2 & p_3 & p_4 \end{pmatrix}$

Dacă $d = 2$, atunci $V = (V_1, V_2)$ cu $P(V_1 = i, V_2 = j) = \frac{p_i p_j}{1 - p_i}$

Dacă $d = 3$, atunci $V = (V_1, V_2, V_3)$ cu $P(V_1 = i_1, V_2 = i_2, V_3 = i_3) = \frac{p_{i_1} p_{i_2} p_{i_3}}{(1 - p_{i_1})(1 - p_{i_1} - p_{i_2})}$ etc.

Un caz particular: $X \sim \begin{pmatrix} 1 & 2 & 3 & 4 \\ .1 & .2 & .3 & .4 \end{pmatrix}$. Dacă $d = 2$ avem 12 de aranjamente, cu următoarele probabilități

12	13	14	21	23	24	31	32	34	41	42	43
2/90	3/90	4/90	2/80	6/80	8/80	3/70	6/70	12/70	4/60	8/60	12/60
0.0222	0.0333	0.044	0.025	0.75	0.1	0.042	0.0857	0.1714	0.0667	0.1333	0.2

Ca să ne convingem că așa stau lucrurile, facem 10000 de simulări

```
> for (j in 1:nsim)
+ {xsim=sample.int(4,2,prob=p); vax[j]=10*xsim[1]+xsim[2]}
> table(vax)
vax
 12   13   14   21   23   24   31   32   34   41   42   43
207  307  439  272  729  948  414  847 1753  685 1319 2080
```

În concluzie, în cazul variabilelor aleatoare unidimensionale discrete (cu o mulțime finită de valori) problema simulării este rezolvată complet prin instrucțiunea `sample`.

Vom prezenta algoritmul general de simulare al unei repartiții care nu face parte din cele furnizate de “R”: acesta este **algoritmul cuantilei**. Ideea este că dacă $U \sim U(0,1)$, atunci

$$X = F^{-1}(U) \text{ este o variabilă aleatoare cu funcția de repartiție } F \quad (6.5.1)$$

Apare o problemă tehnică, aceea că inversa $F^{-1}(p)$ nu există întotdeauna. De aceea o înlocuim cu cuantila de la nivel p , definită prin

$$F^{-1}(p) = x \Leftrightarrow F(x-0) \leq p, F(x) \geq 0 \quad (6.5.2)$$

Metoda: găsim două valori y_1 și y_2 cu proprietatea că $F(x_1) < p$, $F(x_2) > p$, destul de apropiate și aproximăm funcția F cu o funcție de gradul I care trece prin $(x_1, F(x_1))$ și $(x_2, F(x_2))$. Un algoritm simplu și eficient dacă diferența dintre x_1 și x_2 nu este prea mare. În loc să lucrăm cu funcția de repartiție, putem lucra cu coada ei, $\underline{F} = 1 - F$, dar atunci p trebuie înlocuit cu $1 - p$.

Prezentăm două scripturi: `cuantinc` care calculează cuantila plecând de la un tabel $(x, F(x))$ și scriptul `cuantdec`, care o calculează folosind coada ei $(x, \underline{F}(x))$. Aici x este un vector cu valori ale lui x iar y este vectorul cu valorile lui F sau \underline{F} .

```
cuantinc<-function(x,y,p)
{ k=length(which(y<p)); x1=x[k]; x2=x[k+1]
rez=x[k]+(x[k+1]-x[k])*abs((y[k]-p)/(y[k]-y[k+1]))
rez}

cuantdec<-function(x,y,p)
{ p=1-p; k=length(which(y>p)); x1=x[k]; x2=x[k+1]
rez=x[k]+(x[k+1]-x[k])*(y[k]-p)/(y[k]-y[k+1])
rez}
```

Exemplu de aplicare: cuantila unei repartiții Pareto(a, b, c). Coada sa este dată de $\underline{F}(x) = \left(\frac{1}{1+ax^b} \right)^c$, $x \geq 0$. Cuantila se poate calcula și analitic, de data aceasta:

$$\underline{F}(x) = p \Leftrightarrow \underline{F}(x) = 1 - p \Leftrightarrow x = \left(\frac{(1-p)^{-\frac{1}{c}} - 1}{a} \right)^{\frac{1}{b}}. \text{ După cum se vede în scriptul de mai}$$

jos, dacă se ia un pas de .01 pentru x , atunci valoarea exactă x_2 dată de formula de mai sus coincide cu valoarea aproximativă, x_1 , dată de funcția `cuantdec`:

```
fz<-function(x,a,b,c) {y=1/(1+a*x^b)^c;y}
a=1;b=2;c=3
x=seq(0,10,by=.01);y=fz(x,a,b,c);p=.9
x1=cuantdec(x,y,p)
x2=((1-p)^(-1/c)-1)/a^(1/b)
x1 1.074484
x2 1.074446
```

Exemplu în care nu avem nici o formulă analitică, ci doar un tabel de mortalitate. Să se simuleze timpii reziduali de viață pentru un lot de 100 de bărbați de 60 de ani.

Soluție. Folosim tabelul de mortalitate din 2008.

```
## O simulare: N timpi de viata reziduali, barbati de 60
#de ani
b=tabmort(nrbb08,nrbm08);x=0:600/10;y=x
for (j in 0:600) {y[j+1]=tpx(60,x[j+1],b)}
N=100;bb=1:N;for (k in 1:N) {bb[k]=cuantdec(x,y,runif(1))}
> summary(bb)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
```

0.2468 11.3200 18.8100 17.6100 25.4900 32.4600

Simularea ajută la verificare unor ipoteze. De exemplu putem să ne convingem că dacă 1000 de bărbați de 60 de ani se autoasigură cu 6000 de lei pentru ca la momentul decesului (sau după 20 de ani) urmașii lor (sau ei, dacă sunt în viață) să primească 10.000 lei, se creează un fond suficient:

```
> ## 100 de simulări de sume în depozit
> medie=1:100;for (j in 1:100){
+ N;for (k in 1:N) {bb[k]=quantdec(x,y,runif(1))}
+ depozit=1.05^bb;depozit=6000*depozit;medie[j]=mean(depozit)
+ }
> summary(medie)
      Min.      1st Qu.      Median      Mean      3rd Qu.
Max.
 14171     14423     14552     14556     14671     15013
```

Ar putea chiar să parieze că pot primi 14000 lei!

Vectori aleatori și portofolii

Pentru simularea vectorilor aleatori nu există algoritmi generali eficienți.

Totuși, există un algoritm rapid de generare a unui vector aleator repartizat uniform în interiorul mulțimii $\Delta_{k,S} = \{\mathbf{x} \in \mathbb{R}_+^k \mid x_1 + \dots + x_k = S\}$. Este mulțimea vectorilor de probabilitate n – dimensionali, care se mai numesc și **portofolii**.

Motivul este următorul: la bursă există mai multe active financiare (A_1, \dots, A_k) împărțite în acțiuni. După cum se știe, prețul acțiunilor este aleator.

Să presupunem că dacă se investește 1 leu în activul A_j rentabilitatea este o sumă de bani R_j . Mai precis: dacă un jucător la bursă cumpără acțiuni de 1 leu din activul A_j la momentul t_0 , și vrea să le vândă la momentul t_1 , el nu le va vinde tot cu 1 leu, ci cu o sumă X_j , depinzând de cotația de pe piața financiară a activului A_j . Diferența $R_j = X_j - 1$ se numește **rentabilitatea** operațiunii financiare. Sigur că ea depinde de momentele t_0 și t_1 . Din definiție, variabilele aleatoare R_j sunt întotdeauna mai mari sau egale cu 1.

Dacă în loc să cumpere acțiuni în valoare de 1 leu dintr-un singur activ financiar A_j jucătorul de bursă cumpără de x_1 lei acțiuni din activul A_1 , de x_2 lei acțiuni din activul A_2, \dots atunci rentabilitatea operațiunii va fi $R(\mathbf{x}) = x_1 R_1 + x_2 R_2 + \dots + x_k R_k$. Putem întotdeauna să considerăm suma $S = x_1 + \dots + x_k$ ca fiind egală cu 1, deoarece este clar că $R(S\mathbf{x}) = SR(\mathbf{x})$.

Un portofoliu în sens strict este atunci format din **vectorul alocărilor** $\mathbf{x} = (x_1, \dots, x_k) \in \Delta_{k,S}$ al banilor **alocați** fiecărui activ și din **vectorul rentabilităților** $\mathbf{R} = (R_1, \dots, R_k)$. **Valoarea portofoliului** este $R(\mathbf{x}) = x_1 R_1 + x_2 R_2 + \dots + x_k R_k$, media $ER(\mathbf{x})$ se numește **randamentul portofoliului** iar varianța $\text{Var}(R(\mathbf{x}))$ este **riscul portofoliului**. Fie $\mu_j = ER_j$ și $c_{i,j} = \text{Cov}(R_i, R_j) = EX_i X_j - \mu_i \mu_j$. Atunci

$$ER(\mathbf{x}) = x_1 \mu_1 + x_2 \mu_2 + \dots + x_k \mu_k := \mathbf{x}' \boldsymbol{\mu}, \quad \text{Var}(R(\mathbf{x})) = \sum_{i,j} c_{i,j} x_i x_j = \mathbf{x}' \mathbf{C} \mathbf{x} \quad (6.5.3)$$

Am notat cu \mathbf{x}' transpusul lui \mathbf{x} , adică vectorul linie cu aceleași componente.

Problema portofoliului optim este atunci următoarea: **să se găsească cel mai bun portofoliu**. Doar că așa fiind pusă, nu are sens.

Ca două portofolii să fie comparabile, ar trebui să se construiască din aceeași sumă S . Ne-ar trebui un criteriu de optim. Unul din ele este, care are sens, este: **să se găsească portofoliul de randament maxim și risc minim**. Numai că aceasta de obicei nu are soluții.

Alte idei sunt: dintre toate portofoliile cu un randament acceptabil, să se găsească cel de risc minim; sau, dintre cele de un risc acceptabil, să se găsească cele de randament maxim.

O metodă care permite o abordare generală este **principiul utilității medii**¹⁷. O utilitate este o funcție $u: I \rightarrow \mathbb{R}$ continuă strict crescătoare, unde I este un interval de numere reale. Utilitatea medie a unui portofoliu (\mathbf{x}, \mathbf{R}) este $Eu(\mathbf{x}'\mathbf{R})$. Sigur că pentru ca formula să aibă sens trebuie ca $\mathbf{x}'\mathbf{R} \in I \forall \mathbf{x} \in \Delta_{k,S}$. Atunci problema capătă un sens precis

$$\text{Să se găsească } \mathbf{x} \in \Delta_{k,S} \text{ astfel ca } Eu(\mathbf{x}'\mathbf{R}) = \max \quad (6.5.4)$$

O metodă rapidă de rezolvare, cu o precizie acceptabilă a problemei (5.4) este metoda Monte Carlo: se simulează N portofolii uniform repartizate în $\Delta_{k,S}$ și se alege acela pentru care utilitatea medie este maximă.

Algoritm de generare al unui vector repartizat uniform în $\Delta_{k,S}$.

Se simulează k variabile aleatoare independente repartizate $Exp(1)$ și se împarte la suma lor, după care se înmulțește cu S .

Exemplu : $S = 100, k = 5$.

```
> s=100;k=5; x=rexp(k,1); x=x/sum(x)*s;x
[1] 3.154435 20.061981 38.808900 12.395507 25.579176
```

Exemplu de calcul. Rentabilități normal repartizate, utilitate CARA.

Presupunem că vectorul \mathbf{R} este repartizat normal $N(\mu, C)$. În practică se acceptă, deși este clar că nu există variabile aleatoare normale repartizate cu valori în $[-1, \infty)^k$.

Luăm utilitatea $u(t) = -e^{-rt}$. O utilitate ee această formă se numește CARA cu coeficient de aversiune la risc egal cu r .

Se știe că $\mathbf{x}'\mathbf{R} \sim N(\mathbf{x}'\mu, \sqrt{\mathbf{x}'C\mathbf{x}}^2)$, Cum funcția generatoare de momente a unei repartiții $N(\mu, \sigma^2)$ este $m(t) = e^{t\mu + \frac{t^2\sigma^2}{2}}$ avem că $Eu(\mathbf{x}'\mathbf{R}) = -e^{-r\mathbf{x}'\mu + \frac{r^2\mathbf{x}'C\mathbf{x}}{2}}$. Așadar problema (6.5.4) devine

$$\text{Să se găsească } \mathbf{x} \in \Delta_{k,S} \text{ astfel ca } -e^{-r\mathbf{x}'\mu + \frac{r^2\mathbf{x}'C\mathbf{x}}{2}} = \max \Leftrightarrow \mathbf{x}'\mu - \frac{r}{2}\mathbf{x}'C\mathbf{x} = \max \quad (6.5.5)$$

Continuăm exemplul în care $S = 100$. Fie $\mu = (0, .1, .2, .2, .3)$ și $C = AA'$, unde alegem A o matrice aleatoare 5×5 .

```
> a=floor(10*(runif(25)-.5));dim(a)=c(5,5);C=a%*%t(a);
C=C/10;C
      [,1] [,2] [,3] [,4] [,5]
[1,]  3.9  0.2 -2.4  0.6 -1.8
```

¹⁷ Vezi, de exemplu Gheorghița Zbăganu, *Metode matematice în teoria riscului și actuariat*, Editura Universității București, 2004


```

[2,] 0.2 4.7 1.6 -0.4 2.3
[3,] -2.4 1.6 2.3 -0.6 2.3
[4,] 0.6 -0.4 -0.6 0.8 -0.6
[5,] -1.8 2.3 2.3 -0.6 3.1
> eigen(C)
$values
[1] 8.79649765 4.63101195 0.66399316 0.64015618 0.06834105

```

Instrucțiunea `eigen(C)` am pus-o ca să ne convingem ca C este o matrice pozitiv definită. Coeficienții de corelație între variabile sunt

```

> for (i in 1:5) {for (j in 1:5)
{ro[i,j]=C[i,j]/sqrt(C[i,i]*C[j,j])}}
> ro
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 1.00000000 0.04671418 -0.8013367 0.3396831 -0.5176776
[2,] 0.04671418 1.00000000 0.4866393 -0.2062842 0.6025569
[3,] -0.80133668 0.48663925 1.0000000 -0.4423259 0.8613568
[4,] 0.33968311 -0.20628425 -0.4423259 1.0000000 -0.3810004
[5,] -0.51767758 0.60255689 0.8613568 -0.3810004 1.0000000

```

Funcția `portopt` calculează portofoliul optim prin metoda Monte Carlo. Parametrii sunt

- S = suma investită ($S = 100$)
- r = coeficientul de aversiune la risc. Dacă $r = 0$, brokerul este neutru la risc și portofoliul va consta în a pune toți banii pe activul de randament maxim. Dacă r este mare, banii se vor alocă în așa fel încât riscul să fie minim.
- μ este media lui R : un vector cu componente mai mari ca -1
- C este matricea de covarianță
- N este numărul de simulări

Maximul funcției $m'x - rx'Cx/2$ cu r coeficientul de aversiune la risc

```

portopt<-function(S,r,mu,C,N)
{ k=length(mu);xx=matrix(nrow=N,ncol=k);val=c(rep(0,N))
  for (j in 1:N)
    {x=rexp(k,1);x=S*x/sum(x);xx[j,]=x
      val[j]=mu%*%x-r*x%*%C%*%x/2
    }
  valmax=max(val);j=which(val==valmax)
  x = xx[j,];x}

```

Luăm

```
mu =(0.35 0.35 0.20 0.20 0.30)
```

și C matricea de la 6.5.6. Dispersiile $\text{Var}(R_i)$ sunt elementele de pe diagonală ale lui C :

```
sig = (3.9 4.7 2.3 0.8 3.1)
```

Am ales numerele în așa fel încât activele cele mai rentabile să fie și cele mai riscante iar cel mai sigur să fie R_4 . Activul R_5 are și randament acceptabil și risc acceptabil: ar trebui luat în calcul dacă există aversiune la risc.

```

r=0;xoptim=portopt(S,r,mu,C,N);xoptim
x = 70.61 26.60 0.66 1.17 0.96

```

Dacă făceam calculul exact, ne așteptam ca toți cei 100 de lei să fie investiți în A_1 și A_2 . Algoritmul Monte Carlo ne-a spus să punem $70.61 + 26.60 = 97.21$ lei în A_1 și A_2 .

Să luăm acum un coeficient moderat de risc,

```
r=.1; xoptim = portopt(S,r,mu,C,N); xoptim
      xoptim = 31.07  0.38 44.04 24.46  0.06
```

Deja se vede cum activul A_2 este evitat: are același randament ca A_1 dar este mai riscant.

Mai experimentăm

```
r=1    => xoptim = 22.08  0.58 42.83 33.69  0.82
r=100  => xoptim = 27.31  0.10 42.38 28.99  1.21
```

Dacă r este mare, atunci maximul funcției $x'\mu - \frac{r}{2}x'Cx$ se atinge cam în același punct în care $x'Cx$ este minim. Minimul lui $x'Cx$ (1652.5) se atinge în

```
x = 24.8    0    41.82 33.39    0
```

punct care nu diferă mult de cel de la $r = 100$.

Putem calcula și cuantilele valorii $\xi = x'R$ a portofoliului optim. În cazul studiat, când repartiția este normală, știm că $\xi \sim N(x'\mu, \sqrt{x'Cx})$. De exemplu, dacă $r = .1$, probabilitatea de a fi în pierdere cu portofoliul optim alcătuit se poate calcula prin comenzile

```
> r=1;
x=portopt(S,r,mu,C,N); mux=x%*mu; sig=sqrt(x%*C%*x);
pnorm(0,mux,sig)
[1] 0.2820430
> mux = 24.07753, sig = 42.05184
```

Este o probabilitate mare, fiindcă și varianța portofoliului este mare. Dacă înlocuim C cu $C/10$ am fi obținut

```
> r=1;
x=portopt(S,r,mu,C,N); mux=x%*mu; sig=sqrt(x%*(C/10)%*x);
pnorm(0,mux,sig)
p= 0.0390437; mux = 24.62120; sig = 13.72716
```

Dacă schimbăm funcția de utilitate cu alta – de exemplu cu $u(x) = \sqrt{(x+100)}_+$ atunci nu am mai fi putut aplica metoda de mai sus fiindcă nu există formule analitice cu care să putem calcula $E\sqrt{(X+100)}_+$ dacă $X \sim N(\mu, \sigma^2)$. Atunci ar fi trebuit să înlocuim media cu media aritmetică a unui eșantion de un anumit volum dintr-o populație repartizată $X \sim N(\mu, \sigma^2)$. Se poate simula ușor în „R”. Dar dacă vrem să simulăm un vector n_d dimensional, $X \sim N(\mu, C)$, pachetul de bază nu ne ajută. Sigur că există în R și scripturi gata făcute în acest scop, dar nu în pachetul de bază. Programarea este simplă: trebuie scris $X = AY + \mu$, unde $AA' = C$ și $Y = \text{rnorm}(n_d, 0, 1)$ este un vector cu n_d componente independente repartizate $N(0,1)$.

```
## simularea unei repartitii normale N(mu,C)
## N este numarul de simulari,mu este media, un vector cu nd
## componente, C este o matrice de corelatie ndxnd

simnorm<-function(N,mu,C)
{nd=length(mu);nd1=length(C[1,]);nd2=length(C[,1])
if (nd != nd1) {print ("eroare1, dimensiuni gresite");return()}
else if (nd != nd2) {print ("eroare2, dimensiuni
gresite");return()}
Ov=eigen(C);v=Ov[[1]];O = Ov[[2]]
D=O-O; for (j in 1:nd) {D[j,j]=v[j]}
DD=sqrt(D);A=O%%DD%%t(O)
xx=matrix(nrow=N,ncol=nd)
for (j in 1:N) {xx[j,]=A%%rnorm(nd,0,1)+mu}
## fac si proba empirica
mup=mu;CE=C
for (j in 1:nd) {mup[j]=mean(xx[,j])}
for (i in 1:nd) {for (j in 1:nd) {CE[i,j]=cov(xx[,i],xx[,j])}}
rez=list(xx=xx,muemp=mup, covemp=CE)
rez}
```

Scriptul nu numai că simulează N vectori repartizați $N(\mu, C)$, dar le calculează apoi media și covarianța empirică. Din legea numerelor mari, știm că ele nu trebuie să difere prea mult de EX și $Cov(X)$. Iată rezultatul pentru $N = 10000$, $\mu = (0.35 \ 0.35 \ 0.20 \ 0.20 \ 0.30)$ și C matricea de la (5.6).

```
N=10000;xx=simnorm(N,mu,C);muemp=xx[[2]];covemp=xx[[3]]
> muemp = 0.3705928 0.3555581 0.1921526 0.2123407 0.2905003
> covemp
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 3.9503142 0.1655118 -2.4350401 0.6147029 -1.8183625
[2,] 0.1655118 4.6166988 1.5967822 -0.3969151 2.2854434
[3,] -2.4350401 1.5967822 2.3158990 -0.6054019 2.3049773
[4,] 0.6147029 -0.3969151 -0.6054019 0.7873231 -0.6095916
[5,] -1.8183625 2.2854434 2.3049773 -0.6095916 3.0935673
```

Dacă dorim să vedem cum se comportă componentele lui X , putem folosi funcția `summary`:

```
> xsim=xx[[1]]
> for (j in 1:nd) {print(summary(xsim[,j]))}
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-7.5530 -0.9729  0.3854  0.3706  1.6930  7.8400
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-7.8290 -1.0810  0.3550  0.3556  1.7980  9.7950
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-5.9630 -0.8188  0.1742  0.1922  1.1900  5.4900
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-2.8190 -0.3803  0.2119  0.2123  0.8116  3.5680
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-5.9740 -0.8904  0.2826  0.2905  1.4550  6.9090
```

Simularea poate fi folosită atunci când avem de calculat cantități pentru care nu avem o formulă analitică. De exemplu, dacă vrem să calculăm portofoliul optim $x_0 \in \Delta_{nd,S}$ în așa fel încât $E\sqrt{(x_0'R + 10n_d)_+}$ să fie maximă, va trebui să simulăm un număr de N valori ale vectorului rentabilităților, R și un număr de N_X valori de portofolii uniform repartizate din care să îl alegem cel pentru care media empirică este maximă.

Capitolul 7

Aplicații rezolvate prin funcții predefinite în R

7.1 Teste statistice

Reamintim că orice presupunere privind repartiția necunoscută a unei variabile aleatoare poartă numele de *ipoteză statistică*, iar metodele de verificare a ipotezelor statistice poartă numele de *teste statistice*.

Ipoteza care se verifică se numește *ipoteza nulă* (notată de obicei cu H_0) și orice altă ipoteză admisibilă se numește *ipoteză alternativă* (notată de obicei cu H_a).

Verificarea ipotezei nule se face pe baza unei selecții (eșantion) de volum n , x_1, \dots, x_n , dintr-o populație statistică.

Testele statistice, după scopul lor, se pot clasifica în:

- teste de *comparare* a unor parametri ai populației, ce verifică ipoteze precum compararea mediilor a două populații, compararea mediilor mai multor populații, compararea dispersiilor, etc;
- teste de *omogenitate* sau de *independență* care verifică ipoteze de tipul dependenței sau independenței unor factori de clasificare;
- testele de *concordanță*, care verifică dacă distribuția selecției este conformă cu o anumită distribuție teoretică așa cum ar fi distribuția normală.

Testele de comparare se împart la rândul lor în două categorii fundamentale:

- teste *parametrice*, pentru care se presupune că populațiile din care provin eșantioanele au distribuții cunoscute cu cel puțin un parametru necunoscut;
- teste *neparametrice*, pentru care nu se face nicio presupunere despre distribuțiile populațiilor din care provin eșantioanele.

Testele parametrice (fie δ un parametru al distribuției populației) pot fi [7]:

- teste *unilaterale* (direcționale) în care $H_0: \theta = \theta_0$ și $H_1: \theta < \theta_0$ sau $\theta > \theta_0$.
- teste *bilaterale* (nedirecționale) în care $H_0: \theta = \theta_0$ și $H_1: \delta \neq \theta_0$;

La testele parametrice, dacă valoarea parametrului care apare în ipoteza alternativă este unică, atunci ea se numește ipoteză simplă, altfel se numește ipoteză compusă.

Zona de acceptare a unei ipoteze, numită și *interval de încredere*, este un interval în care se acceptă, printr-un test statistic, ipoteza nulă, căreia i se asociază probabilitatea $1-\alpha$. *Zona de respingere* este intervalul dintr-o distribuție de selecție a unei statistici considerate, în care se respinge ipoteza nulă, căreia i se asociază o probabilitate α . Probabilitatea α este numită *prag de semnificație a testului*.

Se pot produce erori de acceptare sau de respingere pe nedrept a unei ipoteze, numite erori de primă speță sau erori de tipul I și de-a doua speță sau erori de tipul II. *Eroarea de tipul I* este eroarea în care se respinge pe ipoteza nulă când în realitate ea este adevărată, probabilitatea asociată fiind α . Valoarea α se alege de obicei ca având valoarea 0,05. *Eroarea de tipul II* este cea în care se acceptă ipoteza nulă atunci când ea este de fapt falsă, probabilitatea asociată fiind β . *Regiunea critică* (de respingere) a testului reprezintă valorile numerice ale testului statistic pentru care ipoteza nulă va fi respinsă. Nivelul de încredere al unui test este $1-\alpha$, iar puterea testului este $1-\beta$.

Un test statistic, în general, parcurge următorii pași [7]:

1. Se formulază ipotezele statistice: o ipoteză nulă H_0 și o ipoteză alternativă H_1 .
2. Pe baza selecției se calculează o statistică numită statistica testului.
3. Se alege un prag de semnificație pentru testul statistic.
4. Se compară valoarea actuală a statisticii testului cu valoarea teoretică.
5. Se stabilesc regulile de decizie de acceptare sau respingere a ipotezei nule.

Există funcții R predefinite, precum **z.test** care se găsește în pachetul de programe R **TechingDemos**. Acest pachet de funcții R se poate găsi ca arhivă zip la adresa web

<http://cran.r-project.org/web/packages/TechingDemos/index.html>

După descărcarea pe calculatorul personal a pachetului, pentru instalare se poate folosi opțiunea „Install package(s) from local zip files...” din meniul *Packages*. După instalarea cu succes, se poate alege opțiunea „Load package...”, din același meniu pentru încărcarea pachetului TechingDemos.

Există funcții R predefinite, precum **sign.test** care se găsește în pachetul de programe **BSDA** (Basic Statistics and Data Analysis). Acest pachet de funcții R se poate găsi ca arhivă zip la adresa web

<http://cran.r-project.org/web/packages/BSDA/index.html>

Se descarcă, se instalează și se încarcă pachetul BSDA ca în cazul pachetului TechingDemos.

Alte funcții referitoare la teste statistice se mai găsesc în pachetul standard „stats” care se instalează automat la instalarea software-ului R.

7.1.1 Teste statistice pentru un singur eșantion

7.1.1.1 Testul Z pentru un singur eșantion

Testul de concordanță Z (sau testul normal) verifică o **ipoteză referitoare la media unei populații repartizate normal, cu dispersia cunoscută**.

Se consideră selecția (eșantionul) X_1, X_2, \dots, X_n de variabile aleatoare independente și identic repartizate, de repartiție $N(m, \sigma^2)$, cu σ cunoscut. Testul Z se aplică pentru eșantioane este de dimensiune $n \geq 30$.

Ipoteza nulă este

$$H_0: m = m_0,$$

iar ipoteza alternativă

$$H_a: m \neq m_0.$$

Statistica Z se calculează după formula

$$Z = \frac{\bar{x} - m_0}{\sigma / \sqrt{n}},$$

unde \bar{x} reprezintă media de selecție a eșantionului:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$$

Ca urmare a teoremei limită centrală, statistica Z are repartiția normală standard $N(0,1)$ pentru n mare. Avem în acest caz, dacă vom alege un risc α , ipotezele și criteriile de acceptare sau respingere conform cu tabelului următor:

H_0	H_a	Regiunea de respingere
$m = m_0$	$m \neq m_0$	$ z > z_{1-\frac{\alpha}{2}}$
$m = m_0$	$m > m_0$	$z > z_{1-\alpha}$
$m = m_0$	$m < m_0$	$z < -z_{1-\alpha}$

Aplicație: Să presupunem că avem 10 voluntari care au făcut un test de inteligență. Media obținută la același test, referitor la întreaga populație este 75. Vrem să verificăm dacă există o diferență statistică semnificativă (cu nivel de semnificație de 5%) între media eșantionului cu valorile 65, 78, 88, 55, 48, 95, 66, 57, 79, 81 la testul de inteligență și media populației, presupunând că dispersia este cunoscută și egală cu 18.

Soluție: Pentru a rezolva această aplicație, se poate folosi testul Z cu un eșantion.

În pachetul de funcții R `TechingDemos` există funcția predefinită `z.test`, cu prototipul:

```
z.test(x, mu = 0, stdev, alternative = c("two.sided", "less",  
"greater"), sd = stdev, conf.level = 0.95)
```

unde

- `x` = vectorul valorilor din eșantionul considerat
- `mu` = media teoretică populației
- `stdev` = dispersia cunoscută a populației
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral)
- `sd` = alternativa la `stdev`
- `conf.level` = intervalul de încredere

În cazul nostru, o secvență de cod R necesară pentru rezolvarea problemei este următoarea:

```
> x = c(65, 78, 88, 55, 48, 95, 66, 57, 79, 81)
> z.test (x,75,18)      #media=75, dispersia=18

One Sample z-test

data:  x
z = -0.6676, n = 10.000, Std. Dev. = 18.000, Std.
Dev. of the sample
mean = 5.692, p-value = 0.5044
alternative hypothesis: true mean is not equal to 75
95 percent confidence interval:
 60.04369 82.35631
sample estimates:
mean of x
 71.2
```

Astfel sunt calculate: valoarea statisticii care este $-0,6676$, intervalul de încredere pentru media vectorului `x`, etc. Acum putem compara valoarea statisticii testului cu valoarea distribuției student cu $10-1=9$ grade de libertate determinată din tabele sau se poate folosi funcția R pentru determinarea acestei valori:

```
qnorm(p, ...)
```

care calculează valoarea teoretică a statisticii t cu pragul de semnificație al testului p . În cazul nostru:

```
> qnorm(0.975)    #Z_(1-alpha/2)
[1] 1.959964
```

Valoarea calculată a statisticii Z este egală cu -0.6676 și în valoare absolută este mai mică decât valoarea teoretică 1.959964 pentru $\alpha = 0.05$. Așadar, se acceptă ipoteza nulă a testului și astfel se decide că media eșantionului este similară cu media întregii populații.

Alternativ, se poate folosi p -valoarea calculată. **p -valoarea unui test statistic** este cea mai mică valoare a nivelului de semnificație pentru care informația extrasă din eșantion este semnificativă (H_0 adevărată se respinge). Cu un prag de semnificație al testului de 5%, se poate aplica regula: dacă valoarea p este mai mare decât 0.05 atunci se acceptă ipoteza nulă H_0 , iar dacă este mai mică decât 0.05 atunci se respinge ipoteza nulă H_0 în favoarea ipotezei alternative H_1 . În cazul nostru, p -valoarea este $0.5044 > 0.05$ și deci se acceptă ipoteza nulă H_0 .

7.1.1.2 Testul t pentru un singur eșantion

Testul t (sau testul Student) pentru un singur eșantion se folosește pentru a verifica o **ipoteză referitoare la media unei populații repartizate normal, cu dispersia de data aceasta necunoscută**.

Se consideră tot o selecție X_1, X_2, \dots, X_n de variabile aleatoare independente și identic repartizate, de repartiție $N(m, \sigma^2)$, cu σ însă necunoscut.

Ipoteza nulă este

$$H_0: m = m_0,$$

iar ipoteza alternativă

$$H_a: m \neq m_0.$$

Dispersia populației din care provine selecția poate fi estimată prin estimatorul deplasat $S^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ sau prin estimatorul nedeplasat $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$, unde \bar{x} reprezintă media de selecție a eșantionului, $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$.

Dacă volumul selecției este mare ($n \geq 30$) atunci $s^2 \approx S^2$ și se poate folosi statistica Z . Dacă însă volumul selecției este mic ($n < 30$), se poate folosi testul Student în care se calculează statistica testului următoare:

$$t = \frac{\bar{x} - m_0}{s / \sqrt{n}}$$

Această statistică a testului t este repartizată Student cu $n-1$ grade de libertate. În rest se procedează ca la testul Z prezentat anterior.

În pachetul de funcții R „stats” există funcția predefinită `t.test`, cu prototipul:

```
t.test(x, y = NULL, alternative = c("two.sided", "less",
"greater"), mu = 0, paired = FALSE, var.equal = FALSE,
conf.level = 0.95)
```

unde

- `x` = vectorul valorilor din primul eșantion
- `y` = al doilea vector (opțional) de valori folosit pentru teste cu două eșantioane
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral)
- `mu` = media teoretică populației
- `paired` = indicator logic pentru un t-test pereche
- `var.equal` = variabilă logică ce indică dacă cele două varianțe sunt egale sau nu
- `conf.level` = intervalul de încredere

Considerând din nou aplicația de la testul anterior, de data aceasta cu dispersie necunoscută, întâi să calculăm statisticele clasice:

```
> mean(x) #media datelor din vectorul x
[1] 71.2
> sd(x)    #deviatia standart a datelor din vect. x
[1] 15.34637
```

Acum, o secvență de cod R care rezolvă aplicația considerată este următoarea:

```
> x = c(65, 78, 88, 55, 48, 95, 66, 57, 79, 81)
> t.test (x, mu=75)

One Sample t-test

data:  x
t = -0.783, df = 9, p-value = 0.4537
alternative hypothesis: true mean is not equal to 75
95 percent confidence interval:
 60.22187 82.17813
sample estimates:
mean of x
 71.2
```

Astfel sunt calculate: valoarea statisticii care este -0.783 , numărul gradelor de libertate care este $10-1=9$, intervalul de încredere pentru media vectorului x , etc. Acum putem compara valoarea statisticii testului cu valoarea distribuției student cu 9 grade de libertate determinată fie din tabele (Tabelul 2 de la Anexa 1), fie se poate folosi funcția R pentru determinarea acestei valori:

```
qt(p, df, ...)
```

care calculează valoarea teoretică a statisticii t cu pragul de semnificație al testului p și numărul gradelor de libertate df . În cazul nostru:

```
> qt(0.975, 9)
[1] 2.262157
```

Așadar, valoarea statisticii t este egală cu -0.783 și în valoare absolută este mai mică decât valoarea critică 2.262157 cu $\alpha = 0.05$ și deci se acceptă ipoteza testului și se decide că media eșantionului este semnificativ similară cu media întregii populații.

Alternativ, se poate folosi p -valoarea calculată. Cu un prag de semnificație al testului de 5%, se poate aplica regula: dacă valoarea p este mai mare decât 0.05 atunci putem accepta ipoteza nulă H_0 , iar dacă este mai mică decât 0.05 atunci respingem ipoteza nulă H_0 în favoarea ipotezei alternative H_1 . În cazul nostru, p -valoarea este $0.4537 > 0.05$ și deci se acceptă ipoteza H_0 .

7.1.1.3 Testul binomial

Testul neparametric binomial poate fi folosit când avem o variabilă cu două valori. Acesta se bazează pe următoarea formulă:

$$P(k \text{ succese în } n \text{ încercări}) = C_n^k p^k q^{n-k}$$

Adică, probabilitatea de a obține k evenimente din n încercări se calculează ca produs al combinațiilor de n luate câte k multiplicat cu probabilitatea teoretică a succesului ridicată la numărul de succese înmulțit cu probabilitatea teoretică a eșecului ridicată la puterea numărului de eșecuri.

Testul se realizează cu ajutorul statisticii Z :

$$Z = \frac{\bar{x} - p}{\sqrt{pq/n}} \sim N(0,1)$$

În R, în pachetul de programe standard “stats” există funcția predefinită

```
binom.test(x, n, p = 0.5, alternative = c("two.sided",  
"less", "greater"), conf.level = 0.95)
```

unde

- x = numărul de succese sau un vector cu 2 componente: (număr de succese, număr de eșecuri)
- n = numărul de încercări
- p = probabilitatea de succes
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral)
- `conf.level` = nivelul de încredere al testului.

Aplicație: Să se testeze ipoteza [8] că senatul de la o universitate americană spune adevărul sau nu și anume că nu se face discriminare sexuală la admitere. Se aleg 500 de studenți și se găsesc 267 băieți. Se consideră cu nivelul de semnificație 5%.

Soluție: O secvență de cod R ce poate soluționa aplicația este următoarea:

```
> binom.test(267, 500)  
  
Exact binomial test  
  
data: 267 and 500  
number of successes = 267, number of trials = 500,  
p-value = 0.1399  
alternative hypothesis: true probability of success  
is not equal to 0.5  
95 percent confidence interval:  
0.4891844 0.5784114  
sample estimates:  
probability of success  
0.534
```

S-a obținut o p -valoare apropiată de 0.1399, fiind astfel mai mare decât pragul de semnificație 0.05 și astfel se acceptă ipoteza nulă. Prin urmare, este foarte probabil ca senatul să spună adevărul.

7.1.1.4 Testul χ^2

Unul dintre cele mai importante teste statistice pentru cea mai bună potrivire a modelului în cazul unei variabile nominale cu două sau mai multe valori. Ipoteza nulă în acest test este aceea că numărul observațiilor în fiecare categorie este egal cu cel prezis de teorie, iar ipoteza alternativă este aceea că numărul observat este diferit de cel așteptat. Statistica testului se calculează în felul următor: se ia numărul observat (n_i), se scade numărul așteptat (n'_i) și se ridică diferența la pătrat. Cu cât deviația față de ipoteza nulă este mai mare, cu atât mai mare ne așteptăm ca diferența dintre numărul observat și cel așteptat. Prin ridicarea la pătrat, toate aceste diferențe devin pozitive. Apoi, fiecare dintre aceste diferențe se împarte la numărul așteptat și aceste diferențe standardizate:

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i}$$

Ca în majoritatea testelor statistice, cu cât mai mare este diferența dintre valoarea observată și cea așteptată, cu atât mai mare devine valoarea testului statistic. Distribuția testului statistic sub ipoteza nulă se potrivește cu distribuția χ^2 . K. Pearson a arătat că, în cazul în care probabilități $p_i = n'_i / n_i$ nu sunt apropiate de 0 sau 1, iar produsele $n'_i = n_i \cdot p_i$, unde $p_i = f(x_i)$, după estimarea parametrilor, nu sunt prea mici (practic nu sunt mai mici decât 5), funcția considerată are repartiția χ^2 cu $(s-1)-k$ grade de libertate, s fiind numărul de valori observate, iar k numărul parametrilor estimați. Dacă între repartiția de selecție și repartiția teoretică există concordanță, atunci statistica χ^2 definită în relația $\chi^2 = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i}$ trebuie să fie

mai mică și nu va depăși o valoare determinată $\chi^2_{(s-1)-k;\alpha}$ corespunzătoare numărului gradelor de libertate $(s-1)-k$ și pragului de semnificație α dat. Regiunea critică a testului va fi dată de inegalitatea $\chi^2 > \chi^2_{(s-1)-k;\alpha}$ și deci, dacă $\chi^2 \leq \chi^2_{(s-1)-k;\alpha}$ acceptăm ipoteza H_0 , în caz contrar o respingem.

În pachetul standard „stats” există funcția R predefinită

```
chisq.test(x, y = NULL, correct = TRUE,  
           p = rep(1/length(x), length(x)), rescale.p = FALSE,  
           simulate.p.value = FALSE, B = 2000)
```

unde

- x = vectorul valorilor din primul eșantion
- y = vector (opțional) de valori din al doilea eșantion
- `correct` = indicator logic referitor la faptul că se aplică sau nu corecție de continuitate când se aplică testul pentru tabele 2 pe 2 – o jumătate se scade din celelalte
- p = vectorul probabilităților de aceeași lungime ca vectorul x cu valori implicite egale
- `rescale.p` = indicator logic cu valoarea TRUE în cazul în care este necesară rescalarea la suma 1, FALSE în cazul în care suma valorilor lui p nu este 1
- `simulate.p.value` = indicator logic ce specifică dacă se calculează valorile p cu simulare Monte-Carlo
- `B` = întreg ce specifică numărul de duplicate folosite în testul Monte-Carlo.

Aplicație: Considerăm că la o facultate din capitală 70% dintre studenți sunt din București, 10% din Pitești, 10% din Ploiești și 10% din alte județe. Vrem să testăm dacă proporțiile observate în eșantionul considerat diferă semnificativ de cele din ipoteza statistică.

Soluție: Presupunem că frecvențele observațiilor sunt următoarele: 145 studenți din București, 24 din Pitești, 20 din Ploiești și 11 din alte județe.

Următoarea secvență de cod R aplică testul χ^2 pentru aplicația considerată

```
> frecvente = c(145,24,20,11)
> proportii = c(0.7,0.1,0.1,0.1)
> chisq.test(frecvente,p=proportii)
```

Chi-squared test for given probabilities

```
data:  frecvente
X-squared = 5.0286, df = 3, p-value = 0.1697
```

Cuantila χ^2 cu $4-1=3$ grade de libertate se pot determina fi din tabele (Tabelul 3 de la anexa 1), fie în R astfel:

```
> qchisq(0.95,3)
[1] 7.814728
```

Aceste rezultate (p -valoarea este $0.1697 > 0.05$ sau valoarea calculată a statisticii testului este $5.0286 < \text{valoarea cuantilei cu } 3 \text{ grade de libertate} = 7.814728$) arată că proporțiile studenților din eșantionul considerat nu diferă semnificativ de valorile furnizate în ipoteza statistică.

7.1.2 Teste statistice pentru două eșantioane

7.1.2.1 Testul t pentru două eșantioane nepereche

Multe cazuri de analiză statistică implică o comparație între mediile a două colectivități generale. Spre exemplu, un patron al unui restaurant dorește să vadă dacă există diferențe între vânzările realizate înainte și după o campanie de publicitate, un grup de consumatori dorește să vadă dacă există o diferență semnificativă între consumul electric pentru două tipuri de cuptoare cu microunde etc.

Testul t pentru două eșantioane independente (nepereche) este probabil cel mai folosit test statistic și cu siguranță cel mai cunoscut. Utilitatea testului constă în faptul că statisticienii examinează cel foarte des natura a două variabile pentru a afla dacă variabilele sunt asociate sau nu. Testul este folosit ca o metodă de **evaluare a diferențelor mediilor a două grupuri**, care pot fi independente (ca de exemplu, [7] tensiunea arterială a pacienților sub tratament cu un anumit medicament și tensiunea arterială a unui grup de pacienți care primesc placebo) sau dependente (ca de exemplu, tensiunea arterilă a unui grup de pacienți înainte și după ce primesc un anumit medicament). Testul poate fi folosit și pentru eșantioane de dimensiune mică, dar care sunt distribuite aproximativ normal.

În general, testul t este utilizat în următoarele trei situații, diferențiate de situația existentă între dispersiile populațiilor și independența eșantioanelor:

- eșantioane independente, dispersii egale;
- eșantioane independente, dispersii diferite;
- eșantioane dependente (perechi, corelate).

Cel mai des întâlnit caz este acel al unui test t pentru două eșantioane independente, unde componentele din cele două eșantioane nu sunt asociate. Presupunând că se cunosc două eșantioane (x_1, \dots, x_n) și (y_1, \dots, y_m) , $n > 30$, $m > 30$, cu mediile m_1 , respectiv m_2 și că diferența dintre cele două este distribuită normal, ipotezele statistice sunt următoarele:

$$H_0: m_1 = m_2$$

$$H_a: m_1 \neq m_2$$

Statistica testului care se calculează este următoarea:

$$t = \frac{\bar{x} - \bar{y}}{\left(\frac{1}{n} + \frac{1}{m}\right) \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{j=1}^m (y_j - \bar{y})^2}{n + m - 2}}}$$

Aceasta este repartizată Student cu $n+m-2$ grade de libertate.

Regiunea critică este dată de:

$$|t| > t_{\alpha/2, n+m-2}$$

Reamintim că în pachetul standard „stats” se găsește funcția `R.t.test` pe care am folosit-o pentru a aplica testul t pentru un singur eșantion. Prototipul acesteia era

```
t.test(x, y = NULL, alternative = c("two.sided", "less", "greater"),  
mu = 0, paired = FALSE, var.equal = FALSE, conf.level = 0.95)
```

unde

- `x` = vectorul valorilor din primul eșantion
- `y` = al doilea vector (opțional) de valori din al doilea eșantion
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral)
- `mu` = media teoretică populației
- `paired` = indicator logic pentru un t-test pereche
- `var.equal` = variabilă logică ce indică dacă cele două varianțe sunt egale sau nu
- `conf.level` = intervalul de încredere

De data aceasta o să aplicăm funcția `t.test` pentru două grupuri.

Exemplu: Considerăm doi vectori independenți (nepereche) cu câte 13, respectiv 8 componente. Aplicăm testul t pentru a vedea dacă mediile lor sunt egale statistic. Dacă nu specificăm în prealabil, se consideră că dispersiile celor două eșantioane nu sunt egale.

Soluție: Următoarea secvență de cod R aplică testul t pentru aplicația considerată

```
> x <- c(79.98, 80.04, 80.02, 80.04, 80.03, 80.03,
+ 80.04, 79.97, 80.05, 80.03, 80.02, 80.00, 80.02)
> y <- c(80.02, 79.94, 79.98, 79.97, 79.97, 80.03,
+ 79.95, 79.97)

> t.test(x, y)

Welch Two Sample t-test
data:  x and y
t = 3.2499, df = 12.027, p-value = 0.006939
alternative hypothesis: true difference in means is
not equal to 0
95 percent confidence interval:
 0.01385526 0.07018320
sample estimates:
mean of x mean of y
 80.02077  79.97875
```

Se observă că p -valoarea este 0.006939 și nu este mai mare decât 0.05, deci există o diferență statistică între mediile celor două grupuri.

Dacă vrem să considerăm cele două dispersii egale, procedăm astfel:

```
> x <- c(79.98, 80.04, 80.02, 80.04, 80.03, 80.03,
+ 80.04, 79.97, 80.05, 80.03, 80.02, 80.00, 80.02)
> y <- c(80.02, 79.94, 79.98, 79.97, 79.97, 80.03,
+ 79.95, 79.97)

> t.test(x, y, var.equal=TRUE)

Two Sample t-test
data:  x and y
t = 3.4722, df = 19, p-value = 0.002551
alternative hypothesis: true difference in means is
not equal to 0
95 percent confidence interval:
 0.01669058 0.06734788
sample estimates:
mean of x mean of y
```

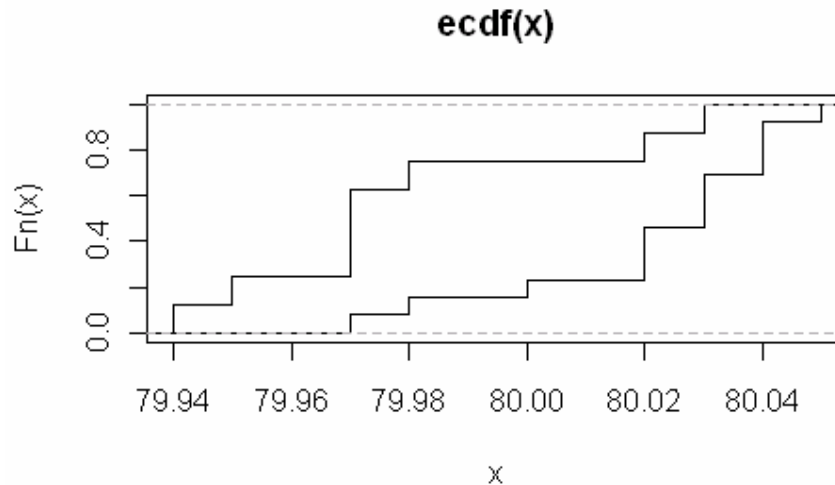
80.02077 79.97875

Rezultatele indică faptul că și de data aceasta (valoarea probabilității p este 0.002551 < 0.05) există o diferență semnificativă statistic între mediile celor două grupuri. Cu alte cuvinte, cum media primului grup este 80.02077 și media celui de-al doilea grup este 79.97875, primul grup este mai semnificativ statistic decât al doilea grup.

Există și modalitatea grafică pentru compararea celor două grupuri de date:

```
>plot(ecdf(x),do.points=FALSE,verticals=TRUE,xlim=range(x, y))  
>plot(ecdf(y),do.points=FALSE,verticals=TRUE,add=TRUE)
```

care conduce la următorul grafic al funcțiilor de repartiție empirice



7.1.2.2 Testul Wilcoxon pentru două eșantioane nepereche

Testul Wilcoxon (al rangurilor cu semn) pentru două eșantioane nepereche este o variantă neparametrică a testului t prezentat anterior aplicat pentru două eșantioane independente (nepereche). Doar că testul Wilcoxon nu face presupunerea că diferența dintre cele două variabile este distribuită normal.

Presupunerea care se face este aceea că mediana diferențelor între perechi de observații este zero. La testul t se făcea presupunerea că media diferențelor perechilor de observații este zero.

Ideea acestui test este următoarea. În primul rând, valorile perechi pentru care diferența este zero se ignoră. Apoi, fiecare valoare absolută a diferențelor dintre observații primește un rang: cea mai mică diferență primește valoarea 1, următoarea rangul 2, etc. Dacă avem diferențe egale, ambele vor primi ca rang media rangurilor. Rangurile tuturor diferențelor dintr-o direcție (pozitiv/negativ) se însumează și rangurile din direcția opusă de asemenea și se obțin valorile W și W_+ . Cea mai mică dintre aceste valori este valoarea testului, W . Spre deosebire de alte teste, valori mici pentru W sunt puțin probabile în cazul ipotezei nule.

În pachetul de funcții R „stats” există funcția predefinită `wilcox.test`, cu prototipul:

```
wilcox.test(x, y = NULL, alternative = c("two.sided", "less",  
"greater"), mu = 0, paired = FALSE, exact = NULL, correct = TRUE,  
conf.int = FALSE, conf.level = 0.95)
```

sau

```
wilcox.test(formula, data, subset, na.action)
```

unde

- x = vectorul valorilor din primul eșantion
- y = al doilea vector (opțional) de valori folosit pentru teste cu două eșantioane
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral)
- μ = media teoretică populației cu valoarea implicită 0
- `paired` = indicator logic pentru un test pereche, implicit cu valoarea FALSE
- `exact` = indicator logic ce specifică dacă ar trebui calculată o p -valoare exactă
- `correct` = indicator logic ce specifică dacă ar trebui aplicată corecție de continuitate la aproximarea normală a valorii lui p
- `conf.int` = indicator logic ce specifică dacă ar trebui calculat un interval de încredere
- `conf.level` = intervalul de încredere
- `formula` = este o formulă de forma membrul_stâng ~ membrul_drept, unde membrul_stâng este o variabilă numerică ce dă valorile, iar membrul_drept este un factor cu două nivele dându-se corespunzătoare celor două grupuri
- `data` = o matrice opțională sau data frame ce conține variabilele din formula `formula`, implicit acestea fiind luate din `environment`;
- `na.action` = o funcție care indică ce ar trebui să se întâmple în cazul în care datele conțin valori NA.

Aplicație: Să considerăm o companie care produce baterii, printre care și baterii scumpe. Însă se dorește modificarea necostisitoare a celor scumpe astfel încât să se ajungă la creșterea duratei de viață. Ipotezele statistice ar fi următoarele:

H_0 : durata de viață a bateriilor modificate este aceeași cu cea a celor scumpe

H_a : durata de viață a bateriilor modificate este mai mare cu cea a celor scumpe

Presupunem că avem la dispoziție un eșantion de 6 baterii scumpe și un eșantion de 5 baterii modificate. Cele 11 sunt construite în condiții identice și independent. Datele disponibile se referă la durata de viață a bateriilor considerate și sunt următoarele: $X = (48, 53, 74, 111, 113, 335)$ și $Y = (62, 101, 167, 174, 190)$.

Următoarea secvență de cod R aplică testul Wilcoxon pentru exemplul considerat:

```
> x=c(48,53,74,111,113,335)
> y=c(62,101,167,174,190)
> wilcox.test(x,y)

      Wilcoxon rank sum test
data:  x and y
W = 10, p-value = 0.4286
alternative hypothesis: true location shift is not
equal to 0
```

Cum p -valoarea este $0.4286 > 0.05$ rezultă că diferența este nesemnificativă statistic între cele două eșantioane.

7.1.2.3 Testul t pentru două eșantioane pereche

Testul t pentru două eșantioane dependente (pereche) este folosit, în general, când măsurătorile provin de la același subiect înainte și după o anumită acțiune. Ipoteza nulă pe care o testăm este aceea că diferența medie între perechile observate este zero. Pentru a aplica acest test mai trebuie ca cele două eșantioane să aibă aceeași lungime.

Presupunem că se cunosc două eșantioane dependente de lungime egală: (x_1, \dots, x_n) și (y_1, \dots, y_n) . În acest caz, statistica testului este următoarea:

$$t = \frac{\bar{x} - \bar{y}}{d},$$

unde \bar{x} , respectiv \bar{y} reprezintă media de selecție celor două eșantioane, iar d reprezintă eroarea standard a diferenței dintre componentele celor două eșantioane care se poate calcula astfel

$$d = \sqrt{\frac{\sum_{i=1}^n (x_i - y_i)^2 - \frac{1}{n} \left(\sum_{i=1}^n (x_i - y_i) \right)^2}{n-1}}$$

După calcularea lui t dependent va trebui să comparăm valoarea obținută cu valoarea t teoretică.

Aplicație: Un cercetător dorește să studieze influența metodelor de relaxare asupra reducerii conduitei agresive la adolescenți. O grupă de 12 subiecți cu tulburări comportamentale este testată inițial în ce privește nivelul agresivității. După acest moment, cei 12 adolescenți sunt învățați să practice relaxarea și urmează sedințe de relaxare de câte 1 oră de 5 ori pe săptămână. După o perioadă de 2 luni, subiecții sunt retestați utilizând aceeași probă pentru a observa dacă nivelul lor de agresivitate a scăzut în urma practicării metodelor de relaxare. Datele obținute de cei 12 adolescenți la chestionarul de agresivitate pe parcursul celor două testări sunt următoarele:

înainte = (21, 18, 15, 11, 24, 23, 17, 19, 19, 17, 20, 15)

după = (15, 13, 16, 13, 13, 12, 11, 10, 15, 17, 14, 16)

Secvența de instrucțiuni R ce rezolvă exemplul considerat este următoarea

```
> inainte = c(21,18,15,11,24,23,17,19,19,17,20,15)
> dupa = c(15,13,16,13,13,12,11,10,15,17,14,16)

> t.test(inainte, dupa, paired=TRUE)

    Paired t-test
data:  pre_test and post_test
t = 3.3726, df = 11, p-value = 0.006223
alternative hypothesis: true difference in means is
not equal to 0
95 percent confidence interval:
 1.56327 7.43673
sample estimates:
mean of the differences
      4.5
```

p -valoarea este 0.006223 și nefiind chiar nulă rezultă că se acceptă ipoteza nulă. Același lucru se poate determina și comparând valoarea testului cu valoarea quantilei distribuției t care se poate afla în R astfel


```
> qt(0.975, 11)
[1] 2.200985
```

Se observă că valoarea obținută 3.3726 este mai mare decât valoarea teoretică 2.200985 și deci mediile celor două eșantioane pereche sunt semnificativ diferite statistic.

7.1.2.4 Testul F pentru două eșantioane independente

Testul F (testul Snedecor) pentru două eșantioane independente este un test de comparare folosit pentru a **determina dacă dispersiile acestora sunt egale sau nu** statistic. Testarea ipotezei privind dispersia poate fi utilizată pentru a trage concluzii privitoare la consistența unor procese economice ori privitoare la riscurile asociate.

Situația când se poate aplica testul F poate fi recunoscută prin:

- două populații caracterizate de variabilele X_1 respectiv X_2 ;
- variabilele sunt repartizate normal, $X_1 \sim N(m_1, \sigma_1^2)$ și $X_2 \sim N(m_2, \sigma_2^2)$;
- două eșantioane, unul din fiecare populație: (x_1, \dots, x_{n_1}) , respectiv (y_1, \dots, y_{n_2}) .

Ipotezele testului pot fi de tip atât de tip lateral cât și de tip bilateral.

Testul bilateral:

$$H_0: \sigma_1^2 / \sigma_2^2 = 1 \quad (\text{egalitatea dispersiilor celor două populații})$$

$$H_a: \sigma_1^2 / \sigma_2^2 \neq 1$$

Testul unilateral:

$$H_0: \sigma_1^2 / \sigma_2^2 = 1$$

$$H_a: \sigma_1^2 / \sigma_2^2 < 1 \text{ sau } \sigma_1^2 / \sigma_2^2 > 1$$

Când ipoteza nulă este adevărată, statistica

$$F^* = \frac{S_1^2}{S_2^2},$$

este repartizată Snedecor F_{α, n_1-1, n_2-1} pentru un prag de semnificație α , unde

$$S_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (x_i - \bar{x})^2 \quad \text{estimație nedeplasată pentru } \sigma_1^2$$

$$S_2^2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (y_i - \bar{y})^2 \quad \text{estimație nedeplasată pentru } \sigma_2^2$$

$$\bar{x} = \frac{1}{n_1} \sum_{i=1}^{n_1} x_i \quad \text{media de selecție a primului eșantion}$$

$$\bar{y} = \frac{1}{n_2} \sum_{i=1}^{n_2} y_i \quad \text{media de selecție a celui de-al doilea eșantion}$$

Repartiția F (Snedecor) cu v_1 , respectiv v_2 grade de libertate are funcția de densitate de forma:

$$f(x) = \left(\frac{\nu_1}{\nu_2}\right)^{\nu_1/2} \frac{\Gamma\left(\frac{\nu_1 + \nu_2}{2}\right)}{\Gamma\left(\frac{\nu_1}{2}\right) \cdot \Gamma\left(\frac{\nu_2}{2}\right)} x^{\frac{\nu_1}{2}-1} \left(1 + \frac{x}{\nu_2}\right)^{-\frac{\nu_1 + \nu_2}{2}}, x \geq 0$$

Decizia, la pragul de semnificație α , pentru testul *bilateral* este următoarea: se respinge ipoteza nulă H_0 în favoarea ipotezei alternative H_a dacă:

$$F^* > F_{1-\alpha/2, n_1-1, n_2-1} \text{ sau } F^* < F_{\alpha/2, n_1-1, n_2-1}$$

Decizia, la pragul de semnificație α , pentru testul *unilateral* este următoarea: se respinge ipoteza nulă H_0 în favoarea ipotezei alternative H_a dacă:

$$F^* > F_{1-\alpha, n_1-1, n_2-1}$$

În pachetul de funcții R „stats” există funcția predefinită `wilcox.test`, cu prototipul:

```
var.test(x, y, ratio = 1, alternative = c("two.sided", "less",  
"greater"), conf.level = 0.95, ...)
```

sau

```
var.test(formula, data, subset, na.action, ...)
```

unde

- `x` = vectorul valorilor din primul eșantion
- `y` = al doilea vector (opțional) de valori folosit pentru teste cu două eșantioane
- `ratio` = raportul presupus pentru dispersiile populațiilor din care provin cele două eșantioane
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral)
- `conf.level` = intervalul de încredere
- `formula` = este o formulă de forma `membrul_stâng ~ membrul_drept`, unde `membrul_stâng` este o variabilă numerică ce dă valorile, iar `membrul_drept` este un factor cu două nivele dându-se corespunzătoare celor două grupuri
- `data` = o matrice opțională sau data frame ce conține variabilele din formula `formula`, implicit acestea fiind luate din `environment`;
- `subset` = vector (opțional) ce specifică submulțimea de observații ce se va folosi
- `na.action` = o funcție care indică ce ar trebui să se întâmple în cazul în care datele conțin valori NA.

Exemplu: Generăm 50 de variabile distribuite $N(0,2)$ și 30 de variabile distribuite $N(1,1)$. Pentru a vedea dacă cele două eșantioane au dispersiile egale sau diferite se poate folosi următorul cod de instrucțiuni R:

```
> x <- rnorm(50, mean = 0, sd = 2)
> y <- rnorm(30, mean = 1, sd = 1)
> var.test(x, y)      # x și y au aceeași dispersie?

      F test to compare two variances
data:  x and y
F = 5.4473, num df = 49, denom df = 29, p-value =
6.182e-06
```

```

alternative hypothesis: true ratio of variances is not
equal to 1
95 percent confidence interval:
 2.736839 10.248600
sample estimates:
ratio of variances
 5.447279

> qf(0.975,49,29) #quantila F(1-alpha/2,n1-1,n2-1)
[1] 1.990354

```

Cum valoarea calculată (5.4473) este mai mare decât valoarea teoretică (1.990354), atunci se respinge ipoteza nulă.

7.1.2.5 Testul Kolmogorov-Smirnov

Se folosește pentru cazul când dorim să determinăm dacă două eșantioane sunt semnificativ diferite. Spre deosebire de alte teste, testul Kolmogorov-Smirnov nu presupune nimic despre distribuția datelor și se poate spune că acest test este neparametric și liber de distribuții.

În general, testul Kolmogorov-Smirnov se poate aplica pentru un eșantion sau pentru două eșantioane. Când se aplică testul pentru un eșantion, se compară distribuția testată cu alte distribuții cunoscute, cum ar fi distribuția uniformă, normală sau log-normală. Testul Kolmogorov-Smirnov pentru două eșantioane compară dacă cele două eșantioane de date provin din aceeași distribuție, deși nu se specifică exact tipul distribuției.

Din studierea convergenței funcției empirice de repartiție F_n către funcția teoretică de repartiție F , Kolmogorov a demonstrat următoarea teoremă:

$$\lim_{n \rightarrow \infty} P\left(d_n < \frac{\lambda}{\sqrt{n}}\right) = K(\lambda) = \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2 \lambda^2},$$

unde $\lambda > 0$ și $d_n = \max |F_n(x) - F(x)|$.

Funcția K este calculată în tabele pentru diverse valori ale lui λ (tabelul distribuției Kolmogorov).

Cu ajutorul acestei teoreme se poate da un criteriu de verificare a ipotezei H_0 că repartiția empirică urmează o anumită lege de repartiție.

Dacă ipoteza H_0 este adevărată, atunci diferențele $|F_n(x) - F(x)|$ nu vor depăși o anumită valoare $d_{\alpha;n}$ pe care o fixăm astfel încât: $P(d_n > d_{\alpha;n} / H_0) = \alpha$, unde α este riscul de gradul întâi. Dar $P(d_n > d_{\alpha;n}) = 1 - P(d_n \leq d_{\alpha;n})$.

Luând $d_{\alpha;n} = \frac{\lambda_\alpha}{\sqrt{n}}$, înseamnă că atunci când H_0 este adevărată și n suficient de mare avem: $P\left(d_n > \frac{\lambda_\alpha}{\sqrt{n}}\right) = 1 - P\left(d_n \leq \frac{\lambda_\alpha}{\sqrt{n}}\right) = 1 - K(\lambda_\alpha) = \alpha$.

Unui prag de semnificație α dat îi corespunde prin relația $K(\lambda_\alpha) = 1 - \alpha$ o valoare λ_α astfel încât, pentru un volum n dat al selecției găsim valoarea $d_{\alpha;n} = \frac{\lambda_\alpha}{\sqrt{n}}$.

Regiunea critică pentru ipoteza H_0 este dată de relația $d_n > \frac{\lambda_\alpha}{\sqrt{n}}$. Deci:

- dacă $d_n < \frac{\lambda_\alpha}{\sqrt{n}}$, există concordanță între F_n și F și se acceptă ipoteza H_0
- dacă $d_n \geq \frac{\lambda_\alpha}{\sqrt{n}}$, nu există concordanță și respingem ipoteza H_0 .

În pachetul standard „stats” din R există funcția `ks.test`

```
ks.test(x, y, alternative = c("two.sided",  
                             "less", "greater"), exact = NULL)
```

unde

- `x` = este un vector de valori din primul eșantion
- `y` = este (opțional) un vector de valori din al doilea eșantion sau un șir de caractere ce reprezintă funcția de distribuție ca de exemplu `pnorm`.
- `alternative` = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral).
- `exact` = NULL sau un indicator logic ce stabilește dacă p -valoarea trebuie calculat; nu se folosește pentru unul sau două eșantioane.

Exemplu: Să considerăm două eșantioane: primul cu 50 de componente repartizate $N(0,1)$, iar al doilea cu 30 de componente repartizate $U(0,1)$.

```
> x <- rnorm(50)  
> y <- runif(30)  
  
> ks.test(x, y) # test Kolmogorov-Smirnov (au  
+aceeași distribuție?)  
  
Two-sample Kolmogorov-Smirnov test  
data: x and y  
D = 0.54, p-value = 1.598e-05  
alternative hypothesis: two-sided
```

Cum p -valoarea este $1.598 \cdot 10^{-5} < 0.5$, rezultă că cele două grupuri sunt semnificativ diferite statistic.

7.1.2.6 Testul Sign

Testul Sign (al semnelor) este unul dintre cele mai simple teste statistice. Acesta se aplică când nu putem măsura diferența dintre cele două eșantioane, dar putem observa că **există o diferență între eșantioanele** în discuție. Se utilizează semnul diferenței și nu valoarea acesteia, atunci când ambele valori sunt măsurate pentru aceiași subiecți. Dacă nu ar exista nicio diferență între valorile perechi, atunci numărul diferențelor pozitive ar trebui să fie egal cu cel al diferențelor negative. Cu cât numărul diferențelor de un anumit semn este mai mare comparativ cu cel al diferențelor de semn opus, cu atât crește posibilitatea ca diferența dintre variabile să fie statistic semnificativă. Putem formula ipoteza testului astfel: $p = P(X > Y)$ și ipoteza nulă este în acest caz

$$H_0: p=0.5$$

Această ipoteză implică faptul că dându-se o pereche aleatoare de măsurători $(x_i, y_i)_i$ atunci x_i și y_i au aceeași șansă de a fi mai mari unul decât celălalt.

Variabilele celor două eșantioane trebuie să fie de tip numeric, iar valorile să fie exprimate în aceeași unitate de măsură, pentru a se putea face diferența lor.

În pachetul de programe BSDA există funcția R predefinită

```
SIGN.test(x, y = NULL, md = 0, alternative = "two.sided",
           conf.level = 0.95)
```

unde

- x = este un vector de valori din primul eșantion (valori NA sau Inf se acceptă, însă sunt ignorate)
- y = este (opțional) un vector de valori din al doilea eșantion (valori NA sau Inf se acceptă, însă sunt ignorate)
- md = valoarea mediane populației specificată în ipoteza nulă
- $alternative$ = direcția ipotezei alternative, valoarea implicită fiind "two.sided" (bilateral).
- $conf.level$ = intervalul de încredere

Exemplu: Să considerăm doi vector ce conțin valori înainte de un eveniment și după acesta. Vom testa ipoteza nulă că „Nu există diferențe între cei doi vectori de măsurători”, cu ipoteza alternativă „Există diferențe între cei doi vectori de măsurători”

```
> x=c(15,14,14,13,15,13,12,13)
> y=c(16,14,15,15,17,15,14,15)

> SIGN.test(x,y)

Dependent-samples Sign-Test
data:  x and y
S = 0, p-value = 0.01563
alternative hypothesis: true median difference is
not equal to 0
95 percent confidence interval:
 -2.000 -0.675
sample estimates:
median of x-y
      -2

               Conf.Level L.E.pt U.E.pt
Lower Achieved CI      0.9297    -2 -1.000
Interpolated CI       0.9500    -2 -0.675
Upper Achieved CI      0.9922    -2  0.000
```

Cum p -valoarea este 0.01563 este foarte mică, atunci se respinge ipoteza nulă în favoarea ipotezei alternative.

7.1.2.7 Testul de independență χ^2

Testul de independență χ^2 se folosește pentru a **verifica dacă două variabile aleatoare sunt independente sau nu**. Se consideră un tabel $s \times k$ în care cele s linii corespund la s selecții independente, iar cele k coloane se referă la k evenimente E_1, \dots, E_k . Se notează cu n_{ij} frecvența observată a evenimentului E_j în selecția de ordin i (numerele n_{ij} se numesc celule de frecvență). Teoria statistică spune că dacă $(n_{i\cdot})_{i=1,2,\dots,s}$ sunt cunoscute apriori, dacă sunt obținute s selecții independente, dacă parametrii $\theta_1, \theta_2, \dots, \theta_l$ sunt estimați cu ajutorul celulelor de frecvență observate și metodei minimului χ^2 , substituind parametrii $\theta_1, \theta_2, \dots, \theta_l$ în $p_j = P(E_j)$ și aceștia în formula

$$X^2 = \sum_{i=1}^s \sum_{j=1}^k \frac{(n_{ij} - n_i p_j)^2}{n_i p_j},$$

cu

$$n_i = \sum_{j=1}^k n_{ij}, i = 1, 2, \dots, s,$$

atunci variabila aleatoare X^2 , pentru n mare, are aproximativ o repartiție χ^2 cu numărul gradelor de libertate egal cu numărul celulelor de frecvență micșorat cu numărul de parametri estimați. În cazul unei selecții de volum n , avem $s = 1$ și cantitatea X^2 rezultată are pentru un număr n mare aproximativ o repartiție χ^2 cu numărul gradelor de libertate egal cu numărul de celule micșorat cu numărul parametrilor estimați și cu o unitate.

Fie p_{ij} probabilitatea ca un individ luat la întâmplare din populația considerată să aparțină liniei j și coloanei j . Dacă $p_{i\cdot}$ este probabilitatea ca un individ să aparțină liniei i și $p_{\cdot j}$ este probabilitatea ca un individ să aparțină coloanei j , atunci ipoteza că cele două variabile sunt independente poate fi scrisă sub forma

$$H_0: p_{ij} = p_{i\cdot} p_{\cdot j}, i = 1, 2, \dots, s \text{ și } j = 1, 2, \dots, k.$$

Dacă avem o selecție de volum n și n_{ij} sunt frecvențele observate, atunci

$$X^2 = \sum_{i=1}^s \sum_{j=1}^k \frac{(n_{ij} - n p_{i\cdot} p_{\cdot j})^2}{n p_{i\cdot} p_{\cdot j}}$$

Deoarece $p_{i\cdot}$ și $p_{\cdot j}$ sunt necunoscute, conform metodei verosimilității maxime, acestea se pot estima cu $n_{i\cdot} / n$, respectiv $n_{\cdot j} / n$, unde

$$n_{i\cdot} = \sum_{j=1}^k n_{ij}, \text{ respectiv } n_{\cdot j} = \sum_{i=1}^s n_{ij}.$$

Atunci variabila

$$X^2 = \sum_{i=1}^s \sum_{j=1}^k \frac{(n_{ij} - n_{i\cdot} n_{\cdot j} / n)^2}{n_{i\cdot} n_{\cdot j} / n}$$

urmează repartiția $\chi^2 [sk-(s+k-2) = (s-1)(k-1)]$.

Testul este valid dacă cel puțin 80% dintre frecvențele probabile depășesc valoarea 5 și toate frecvențele probabile depășesc valoarea 1. Această condiție este necesară pentru ca repartiția multinomială să poată fi aproximată cu o repartiție normală, dar limitează semnificativ utilizarea testului χ^2 . În cazul în care o frecvență probabilă este sub valoarea 2 sau dacă mai mult de 20% din frecvențele probabile sunt sub valoarea 5, se recomandă utilizarea testului exact al lui Fisher. Atenție că frecvențele probabile calculate în cadrul testului nu sunt frecvențe observate.

Aplicație: Să considerăm baza de date „survey”, în care coloana „smoke” conține răspunsurile referitoare la fumat a 237 de studenți, cu variantele „heavy”, „regul” (regular), „ocass” (ocasionally) sau „never”, iar coloana „exer” se referă la cât de multă mișcare fac studenții, cu variantele acceptate „freq” (frecvent), „some” sau „none”. Aceste date se pot pune într-un tabel de contingență astfel:

```
> library(MASS)
> t = table(survey$Smoke, survey$Exer)
> t
```

	Freq	None	Some
Heavy	7	1	3
Never	87	18	84
Occas	12	3	4
Regul	9	1	7

Să testăm acum ipoteza că obiceiul fumatului studenților chestionați este independent de mișcarea fizică făcută de aceștia la nivelul de semnificație de 0,05.

Soluție: Putem aplica testul χ^2 de independență prin utilizarea funcției predefinite `chisq.test` astfel

```
> chisq.test(t)

Pearson's Chi-squared test

data:  t
X-squared = 5.4885, df = 6, p-value = 0.4828

Warning message:
In chisq.test(t) :
  Chi-squared approximation may be incorrect
```

Cum p-valoarea găsită are valoarea 0,4828 care este mai mare decât pragul de semnificație de 0,05, atunci acceptăm ipoteza nulă conform căreia obiceiul fumatului este independentă de nivelul de mișcare fizică a studentului. Mesajul de eroare se datorează numărului de celule din tabelul de contingență. Pentru a evita acest avertisment, putem combina coloanele a doua cu a treia și noul tabel sa-l numim `t2` și apoi reaplicăm testul pentru acest nou tabel.

```
> t2=cbind(t[, "Freq"], t[, "None"]+t[, "Some"])
> t2
```

	[,1]	[,2]
Heavy	7	4
Never	87	102
Occas	12	7
Regul	9	8

```
> chisq.test(t2)
      Pearson's Chi-squared test
data:  t2
X-squared = 3.2328, df = 3, p-value = 0.3571
```

Din nou se observă din nou independența celor două variabile considerate.

7.1.2.8 Testul de independență al lui Fisher

Testul exact al lui Fisher reprezintă o alternativă a testului χ^2 în examinarea asociațiilor în cadrul unui tabel de contingență 2×2 , atunci când frecvențele probabile sunt mici. Condiția de aplicare a acestui test este ca totalurile pe rânduri și pe coloane să fie fixe, cunoscute dinainte.

Pentru a ști dacă între cele două variabile aleatoare componente ale vectorului aleator (X, Y) există o legătură semnificativă, se poate emite ipoteza H_0 că cele două variabile aleatoare sunt independente și că deci repartițiile unităților în cele patru rubrici ale tabelului vor fi următoarele

	x_1	x_2	Total
y_1	$\frac{n_{o1}n_{1o}}{n}$	$\frac{n_{o2}n_{1o}}{n}$	n_{1o}
y_2	$\frac{n_{o1}n_{2o}}{n}$	$\frac{n_{o2}n_{2o}}{n}$	n_{1o}
	n_{o1}	n_{o2}	

Pentru a vedea dac între cele două variabile aleatoare există sau nu o conexiune, trebuie să calculăm nu numai probabilitatea $P(n_{11}, n_{12}, n_{21}, n_{11} / n_{o1} n_{o2})$, dar și probabilitățile de a se obține repartiții care față de cea observată se îndepărtează mai mult de repartiția din cazul independenței.

Dacă $n_{11} < n_{o1}n_{1o} / n$, tabelele ce se îndepărtează mai mult sunt acelea care au prima celulă egală respectiv cu $(n_{11}-1), (n_{11}-2), \dots, 1, 0$ unități.

Dacă $n_{11} > n_{o1}n_{1o} / n$, tabelele extreme sunt acelea care au în prima celulă respectiv $(n_{11}+1), (n_{11}+2), \dots, \min\{n_{o1}, n_{1o}\}$.

Compararea sumelor

$$P(n_{11}) = \sum_{i=0}^{n_{11}} p(i), \text{ dacă } n_{11} < n_{o1}n_{1o} / n$$

$$P(n_{11}) = \sum_{i=n_{11}}^{\min\{n_{o1}, n_{1o}\}} p(i), \text{ dacă } n_{11} > n_{o1}n_{1o} / n$$

În pachetul de programe “stats” există funcția R predefinită

```
fisher.test (x, y = NULL, workspace = 200000, hybrid = FALSE,
  control = list(), or = 1, alternative = "two.sided",
  conf.int = TRUE, conf.level = 0.95,
  simulate.p.value = FALSE, B = 2000)
```


unde

- x = este un tabel de contingență 2×2 sau un prim obiect factor;
- y = este (opțional) un al doilea obiect factor sau ignorat dacă x este matrice;
- $workspace$ = un întreg ce specifică dimensiunea spațiului de lucru folosit în algoritm, cu unitatea de măsură de câte 4 bytes;
- $hybrid$ = o valoare logică implicită FALSE, folosită pentru tabele mai mari de 2×2 , caz în care indică probabilitățile exacte (implicite) sau o aproximare hibridă ce ar trebui calculată;
- $control$ = o listă conținând componente predefinite pentru un control jos al algoritmului; în prezent, singurul folosit este “mult”, un întreg pozitiv ≥ 2 cu valoarea implicită 30, folosit doar pentru tabele mai mari decât 2×2 ; acesta specifică de câte ori ar trebui alocat spațiu pentru căi, cât și pentru valori;
- or = raportul de diferențe din ipoteză (folosit doar pentru tabele 2×2);
- $alternative$ = direcția ipotezei alternative, valoarea implicită fiind “two.sided” (bilateral), dar mai poate fi și “greater” sau “less” – se poate specific doar prima literă din șir; valabil doar pentru tabele 2×2 ;
- $conf.int$ = parametru logic ce indică dacă se dorește determinarea unui interval de încredere intervalul de încredere
- $conf.level$ = nivelul de încredere pentru intervalul de încredere returnat – folosit doar pentru tabele 2×2 dacă $conf.int=TRUE$;
- $simulate.p.value$ = o variabilă logică ce specifică dacă p-valoarea se va calcula prin simulare Monte Carlo – folosit pentru tabele mai mari decât 2×2 ;
- B – un întreg ce specifică numărul de duplicate folosit în testul Monte Carlo.

Pentru a ști dacă între cele două variabile aleatoare componente ale vectorului aleator (X, Y) există o legătură semnificativă, se poate emite ipoteza H_0 că cele două variabile aleatoare sunt independente și deci că repartițiile unităților în cele patru rubrici ale tabelului vor fi următoarele:

	x_1	x_2	total
y_1	$n_{o1}n_{1o} / n$	$n_{o2}n_{1o} / n$	n_{1o}
y_2	$n_{o1}n_{2o} / n$	$n_{o2}n_{2o} / n$	n_{2o}
	n_{o1}	n_{o2}	

Pentru a vedea dacă între cele două variabile aleatoare există sau nu o conexiune, trebuie să calculăm nu numai probabilitatea $P(n_{11}, n_{12}, n_{21}, n_{11} / n_{o1} n_{o2})$, dar și probabilitățile de a se obține repartiții care față de cea observată se îndepărtează mai mult de repartiția din cazul independenței.

Dacă $n_{11} < n_{o1}n_{1o} / n$, tabelele ce se îndepărtează mai mult sunt acelea care au prima celulă egală respectiv cu $(n_{11}-1), (n_{11}-2), \dots, 1, 0$ unități.

Dacă $n_{11} > n_{o1}n_{1o} / n$, tabelele extreme sunt acelea care au în prima celulă respectiv $(n_{11}+1), (n_{11}+2), \dots, \min\{n_{o1}, n_{1o}\}$.

Compararea sumelor

$$P(n_{11}) = \sum_{i=0}^{n_{11}} p(i), \text{ dacă } n_{11} < n_{o1}n_{1o} / n$$

$$P(n_{11}) = \sum_{i=n_{11}}^{\min\{n_{o1}, n_{1o}\}} p(i), \text{ dacă } n_{11} > n_{o1}n_{1o} / n$$

probabilităților relative tabelului observat și cu probabilitățile care se îndepărtează tot mai mult de repartiția în cazul independenței comparată cu nivelul de semnificație α , permite să decidem asupra concordanței sau neconcordanței rezultatelor experimentale cu ipoteza de H_0 de independență.

Într-o asemenea comparație trebuie să ținem seama dacă alternativa ipotezei H_0 este uni- sau bilaterală. În primul caz vom compara pe

$$P(n_{11}, n_{12}, n_{21}, n_{22} | n_{o1}, n_{o2}) = \frac{n_{o1}!n_{o2}!n_{1o}!n_{2o}!}{n!n_{11}!n_{12}!n_{21}!n_{22}!}$$

cu α , într-un al doilea cu $\alpha/2$.

Trebuind să recurgem la factoriale, vom lua în întotdeauna în considerație rubrica în care apare cea mai mică frecvență observată. Aceasta este posibil, deoarece toate tabelele se alctuiesc cu aceeași repartiție marginală și pe de altă parte fiind vorba de tabele dictonice, au un singur grad de libertate, care odată fixată valoarea pentru o celulă, automat rămân fixate și valorile pentru celelalte.

Aplicație: Cluster și Galii în 2002 cu zburat cu un avion pentru a urmări speciile stârcul mare albastru și marele egrete de la locul de odihnă la primul loc de hrănire de pe lacul Peltier, Minnesota și au notat tipul pământ pe care au aterizat, obținându-se următorul tabel

	Stârc	Egretă
Vegetație	15	8
Țărm	20	5
Apă	14	7
Structuri	6	1

Să testăm ipoteza conform căreia cele două specii de păsări folosesc locurile de aterizare în proporții egale.

Soluție:

```
t = matrix (c(15,8,20,5,14,7,6,1), nrow=4, ncol=2,
dimnames = list( Loc = c("vegetatie", "tarm", "apa",
"structuri"), Pasari = c("starc","egreta"))) )
> t
      Pasari
Loc      starc  egreta
vegetatie    15     14
tarm          8      7
apa          20     6
structuri     5      1
> fisher.test(t)

Fisher's Exact Test for Count Data

data:  t
p-value = 0.1587
alternative hypothesis: two.sided
```

Aplicarea testului lui Fisher conduce la o p -valoare de $0.1587 > 0.05$, deci nu este nicio dovadă că cele două specii de păsări folosesc locuri de aterizare în proporții diferite.

7.1.3 Teste statistice pentru 3 sau mai multe eșantioane

7.1.3.1 Testul de analiză dispersională ANOVA

Analiza dispersională (ANOVA, “Analysis of variance”) este cea mai folosită metodă pentru compararea mediilor grupurilor pentru variabilele cantitative. ANOVA combină și extinde testele t și χ^2 , prin testarea egalității dintre trei sau mai multe medii (pentru trei sau mai multe grupuri). De subliniat este faptul că ANOVA se poate folosi cu rezultate foarte bune și la comparația dintre două medii, însă testul își arată adevărata valoare la trei sau mai multe medii.

Se testează așadar legătura dintre o variabilă metrică (pentru care se calculează media) și o variabilă calitativă (a cărei valori sau categorii sunt considerate grupuri independente). De asemenea, ANOVA face o introducere clară în analiza cauzală: variabila cauză (independentă) este cea calitativă, iar variabila efect (dependentă) este cea metrică.

Sunt foarte multe modele experimentale care pot fi analizate cu diferite tipuri de analize dispersionale, însă în această secțiune vom vorbi despre analiza dispersională simplă (într-un singur sens).

Aplicație: Se considerăm următoarele două variabile: “vârsta” = vârsta la angajare a unei persoane (variabilă cantitativă) și “stratin” = strategia de atragere a tinerilor (variabilă calitativă). În cazul nostru, ne raportăm doar la strategia de atragere a tinerilor; este important însă de știut că vârsta la angajare a tinerilor poate fi influențat de diverși factori și că există variante ale analizei de varianță care iau în calcul mai mulți asemenea factori. De pildă, varianta care ia în calcul doi factori se numește ANOVA bi-factorial (în engl. two-way ANOVA), iar varianta care ia în calcul mai mulți factori se numește ANOVA multi-factorial (în engl. multi-way ANOVA).

Să presupunem că studiem oportunitățile de acces a tinerilor pe piața forței de muncă și analizăm diferite strategii folosite pentru a atrage tinerii să se angajeze. Ipoteza pe care dorim să o testăm este următoarea: “media de vârstă a persoanelor nou angajate este influențată de strategia de atragere utilizată”. Setul de ipoteze statistice pentru ANOVA este:

$$H_0 : m_1 = m_2 = \dots = m_k$$

$$H_a : \text{cel puțin două medii sunt diferite}$$

În cuvinte, ipoteza nulă susține că nu este nici o diferență între rezultatele diferitelor strategii (fie sunt toate strategiile foarte bune și atrag mulți tineri, fie sunt toate foarte slabe și nu atrag tineri), iar ipoteza alternativă susține cel puțin o strategie dă rezultate mai bune decât cel puțin una dintre celelalte (este posibil ca o strategie să aibă un rezultat de mijloc, care nu este semnificativ diferit nici față de strategia de succes maxim, nici față de strategia care atrage cei mai puțini tineri; diferența semnificativ în acest caz există doar între prima și ultima strategie).

După cum se poate vedea în tabelul următor, primul grup (cel de control, din localitatea unde nu s-a aplicat nici o strategie) are o medie a vârstei la angajare de 27,8 ani cu o abatere standard de 3,65 ani; al doilea grup (din localitatea unde s-a aplicat prima strategie) are o medie de 23,6 ani cu o abatere standard de 3,34 ani, iar al treilea grup o medie de 26,3 ani cu o abatere standard de 3,59 ani. La o primă

vedere, toate cele trei grupuri conțin tineri: există vreo diferență semnificativă între cele trei medii? Cum testăm, mai exact, acest lucru?

Nr.crt	Localitate 1	Localitate 2	Localitate 3
1	22	28	20
2	27	22	28
3	32	24	31
4	30	18	26
5	29	21	26
6	27	26	30
7	33	25	21
8	24	20	25
9	24	24	29
10	30	28	27
\bar{x}	27,8	23,6	26,3
s	3,65	3,34	3,59

Analizând tabelul, putem extrage câteva informații interesante, care ne vor ajuta în cele ce vor urma. Avem un eșantion total format din 30 de persoane, deci $n=30$. Acest eșantion este format din trei grupuri independente de câte 10 persoane fiecare (subeșantioane din trei localități diferite); avem deci: $n_1=10$, $n_2=10$ și $n_3=10$. Pentru fiecare dintre cele trei localități / grupuri putem calcula câte o medie și câte o abatere standard; mai avem așadar: $\bar{x}_1 = 27,8$ cu $s_1 = 3,65$, $\bar{x}_2 = 23,6$ cu $s_2 = 3,34$ și $\bar{x}_3 = 26,3$ cu $s_3 = 3,59$. În același timp, putem calcula o medie generală pentru eșantionul total (pentru toate cele 30 de observații) $\bar{\bar{x}} = 25,9$ precum și o abatere standard total $s = 14,714$.

Poate fi determinată variația mediilor de grupuri în jurul mediei generale:

$$s_{\bar{x}}^2 = \frac{1}{k-1} \sum_{j=1}^k (\bar{x}_j - \bar{\bar{x}})^2$$

Aceasta indică o estimare a erorii standard din populație, de unde putem extrage foarte simplu varianță din populație $ES = \sigma / \sqrt{n}$.

Ne confruntăm, deci, cu două tipuri de variații: o variație între grupuri (variația mediilor de grup în jurul mediei generale) și una în interiorul grupurilor (variațiile observațiilor în jurul fiecărei medii de grup). Ambele tipuri de variații sunt folosite ca estimări ale variației generale în populație.

Analiza de varianță se bazează pe comparația dintre două estimări ale varianței σ^2 pentru întreaga populație. Logica analizei este următoarea: dacă cele două estimări ale varianței din populație σ^2 sunt aproximativ egale, atunci ipoteza nulă este adevărată (în populație, toate mediile sunt egale). Dacă ipoteza de nul nu este adevărat, atunci cele două estimări ale varianței vor fi semnificativ diferite.

```
> localit1 = c(22,27,32,30,29,27,33,24,24,30)
> localit2 = c(28,22,24,18,21,26,25,20,24,28)
> localit3 = c(20,28,31,26,26,30,21,25,29,27)
> localit = c(localit1,localit2,localit3)
> n=rep(10,3)
> v=rep(1:3,n)
> table(localit)
localit
```

```

18 20 21 22 24 25 26 27 28 29 30 31 32 33
  1  2  2  2  4  2  3  3  3  2  3  1  1  1
>f=function(x) c(mean=mean(x),var=var(x),n=length(x))
> tapply(localitati,v,f)
$`1`
      mean      var      n
27.80000 13.28889 10.00000

$`2`
      mean      var      n
23.60000 11.15556 10.00000

$`3`
      mean      var      n
26.3 12.9 10.0

> date = data.frame(localit=localit,v=factor(v))
> model = lm(localit~v,date)
> anova(model)
Analysis of Variance Table

Response: localit
      Df Sum Sq Mean Sq F value    Pr(>F)
v       2   90.6   45.300   3.6391 0.03987 *
Residuals 27  336.1   12.448
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.'
0.1 ' ' 1

```

Valoarea de 3,639 (mult mai mare decât 1) ne sugerează că ipoteza nulă este pe cale de a fi respinsă, pentru că variația explicată de diferențele dintre grupuri este mult mai mare decât variația datorată erorilor aleatoare; existând diferențe majore între grupuri, vor exista cu siguranță diferențe între mediile acestora. Modalitatea alternativă este de a compara p -valoarea cu pragul de semnificație ales; cum p -valoarea este 0.03987 și este mai mică decât $\alpha = 5\%$ se respinge ipoteza nulă: cel puțin una dintre strategii a dat rezultate.

Anova nu ne spune însă decât dacă acceptăm sau respingem ipoteza nulă, adică dacă există sau nu o medie diferită statistic de celelalte. În cazul acceptării ipotezei nule cu Anova, pentru a vedea unde este diferența, putem folosi un **test al diferențelor dintre medii** și anume **testul Tukey**.

În pachetul standard de funcții R „stats” există funcția predefinită `Tukey.HSD`, cu prototipul:

```
TukeyHSD (x, which, ordered = FALSE, conf.level = 0.95, ...)
```

unde

- `x` = un obiect de potrivire a modelului, de obicei `aov`;
- `which` = vector de caractere ce specifică lista de termeni din modelul de potrivire pentru care ar trebui calculate intervalele;
- `ordered` = valoare logică ce specifică dacă nivelele factorului ar trebui ordonate crescător după medie înainte de a calcula diferențele;
- `conf.level` = nivelul de încredere

Aceasta determină o mulțime de intervale de încredere pentru diferențele dintre mediile nivelelor unui factor. John Tukey a calculat intervale bazate pe valorile mediilor grupurilor și nu pe baza diferențelor individuale.

În cazul nostru, o secvență de cod R necesară pentru rezolvarea problemei este următoarea:

```
> l1 = c(22,27,32,30,29,27,33,24,24,30)
> l2 = c(28,22,24,18,21,26,25,20,24,28)
> l3 = c(20,28,31,26,26,30,21,25,29,27)
> localitati = c(l1,l2,l3)
> n = rep(10,3)
> v = rep(1:3,n)

> grup = factor(v)
#inaintea de functia "aov" trebuie convertit vectorul
#calitativ intr-un factor

> analiza = aov(localitati~grup) #analiza dispersiei
> summary(analiza)
              Df Sum Sq Mean Sq F value    Pr(>F)
grup           2   90.6   45.300    3.6391 0.03987 *
Residuals     27  336.1   12.448
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

> TestulTukey = TukeyHSD(analiza,"grup")
> TestulTukey
      Tukey multiple comparisons of means
    95% family-wise confidence level
Fit: aov(formula = localitati ~ grup)
$grup
      diff      lwr      upr      p adj
2-1  -4.2 -8.112164 -0.2878359 0.0335040
3-1  -1.5 -5.412164  2.4121641 0.6136386
3-2   2.7 -1.212164  6.6121641 0.2194422
```

Se observă că cele mai apropiate medii sunt cele ale grupurilor 1 și 3, apoi ale grupurile 3 și 2. Se observă că *p*-valoarea pentru grupurile 2 și 1 este 0.0335040, valoare mai mică decât 0.05, așadar media grupului 2 este diferită statistic de media grupului 1 la un nivel de semnificație de 5%.

7.1.3.2 Testul Kruskal-Wallis

Acesta este o **variantă neparametrică pentru a testa egalitatea mediilor** între grupuri, identică cu Anova unidirecțională, datele fiind înlocuite de ranguri. În plus, testul Kruskal-Wallis **nu presupune populații distribuite normal**.

Pe baza celor k eșantioane independente din k populații, pașii de calcul ai testului sunt următorii:

- se ordonează valorile din cele k eșantioane ca și cum ele ar fi fost obținute din aceeași populație;
- se atribuie câte un rang observațiilor așezate în ordine, de la 1 la n (pentru observații cu același rang, fiecăruia i se va atribui media rangurilor pe care le-ar fi primit dacă nu ar fi avut același rang).
- se calculează statistica $H_{KW} = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n+1)$, unde R_i reprezintă suma rangurilor din eșantionul i , n_i reprezintă numărul observațiilor eșantionului i , $i=1,2,\dots,k$, iar n este volumul selecției totale;
- se compară H_{KW} cu valoarea teoretică $\chi^2_{(1-\alpha)(k-1)}$ și dacă $H_{KW} \geq \chi^2_{(1-\alpha)(k-1)}$ atunci ipoteza nulă se respinge, respectiv dacă $H_{KW} < \chi^2_{(1-\alpha)(k-1)}$ atunci ipoteza nulă se acceptă.

În pachetul standard de funcții R „stats” există funcția predefinită `kruskal.test`, cu următoarele variante de prototip:

```
kruskal.test(x,...)
sau
kruskal.test(x,g,...)
sau
kruskal.test(formula,data,subset,na.action,...)
```

unde

- x = un vector numeric (eșantionul) sau o listă de vectori numerici (eșantioanele independente);
- g = vector de factori ce ne dau grupurile elementelor corespunzătoare din x ; acest argument este ignorat dacă x este listă;
- `formula` = formulă de forma `valori_date ~ grupurile corespunzătoare`;
- `data` = matrice opțională sau data frame ce conține variabile de forma `formula`.
- `subset` = vector opțional ce specifică o submulțime de observațiile ce pot fi folosite;
- `na.function` = funcție ce specifică ce-ar trebui să se întâmple cu datele ce conțin NA.

Aceasta aplică testul de însumare a rangurilor Kruskal-Wallis pentru ipoteza nulă că parametrii distribuției lui x sunt aceeași în fiecare eșantion, au ipoteza alternativă că ei diferă în cel puțin un parametru.

Aplicăm acum testul Kruskal-Wallis în R pentru aplicația de la Anova unidirecțională:

```
> by(localitati,grup,summary)
grup: 1
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
	22.00	24.75	28.00	27.80	30.00	33.00

grup: 2						
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
	18.00	21.25	24.00	23.60	25.75	28.00

grup: 3						
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
	20.00	25.25	26.50	26.30	28.75	31.00

```

> kruskal.test(localitati,grup)

Kruskal-Wallis rank sum test

data: localitati and grup
Kruskal-Wallis chi-squared = 5.8546, df = 2, p-value
= 0.05354

```

Se observă că p -valoarea este 0.05354, puțin mai mare decât pragul de semnificație 0.05, așadar se acceptă ipoteza nulă conform căreia mediile grupurilor sunt egale statistic.

7.2 Serii de timp

O serie de timp reprezintă o mulțime de date măsurate la intervale egale de timp. Datele sunt numerice și reprezintă valoarea unei mărimi fizice (serie unidimensională) sau a mai multor mărimi măsurate simultan (serie multidimensională sau vectorială). În cele ce urmează o să exemplificăm o parte din teoria seriilor de timp pe date reale. Utilizăm ca date concentrația de CO în centrul orașului București. Datele măsurate la interval de o oră provin de la stația Cercul militar, Calea Victoriei, între 1/1/2007 01:00 și 12/30/2008 24:00. Ca date de lucru am utilizat esantionul dintre 1/3/2008 01:00 - 31/5/2008 24:00 și am încercat să prezicem poluarea pentru următoarele câteva zile.

Datele sunt citite inițial dintr-un fișier text și reprezentate grafic:

```
co=read.table(file.choose())
vco=as.matrix(co)
tsco<-ts(vco,frequency=168)
plot(tsco,main="Concentratia de CO intre 1/3/2008 01:00 -
31/5/200824:00")
```

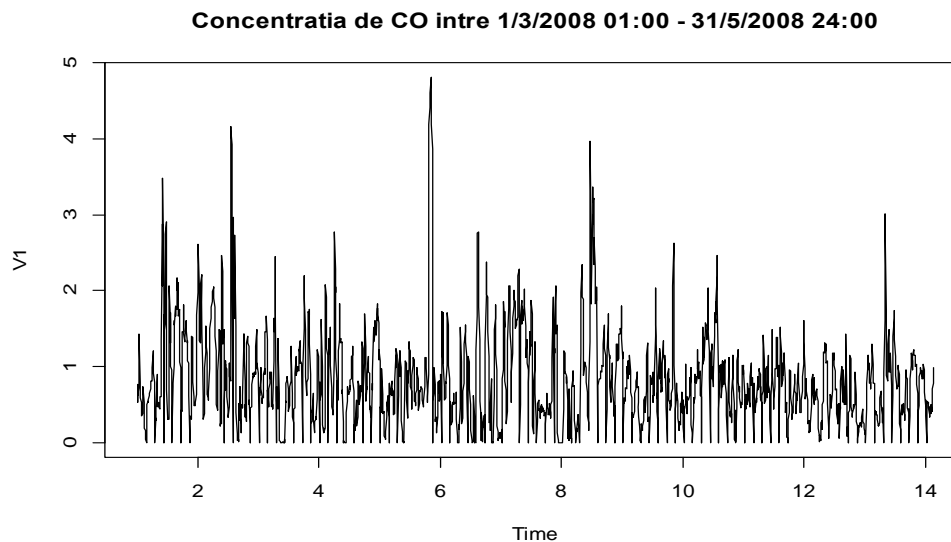


Fig 7.1

Funcțiile din R pentru analiza seriilor de timp prevăd că datele măsurate să fie convertite în obiecte de tip *ts*. Conversia se face prin comenzile *ts(...)* sau *as.ts(...)*, mai exact

```
ts(data = NA, start = 1, end = numeric(0), frequency = 1,
    deltat = 1, ts.eps = getOption("ts.eps"), class = ,
    name= )
as.ts(x, ...)
```

Dintre parametrii din funcțiile de mai sus *data* reprezintă vectorul sau matricea cu datele măsurate, *start* este momentul de începere al observațiilor, *frequency* este

numărul de observații pe unitatea de timp, x este un obiect R. Noi am luat frecvența 168= numărul de ore dintr-o săptămână (bănuind că există o periodicitate săptămânală). Ca urmare unitatea de timp este săptămâna.

7.2.1 Analiza statistică a seriilor de timp

Analiza statistică a seriilor de timp conform cu modelul Box-Jenkins¹⁸ cuprinde următoarele stagii:

- a. *Identificarea modelului*
- b. *Estimarea parametrilor*
- c. *Validarea modelului (diagnosticul veridicității modelului ales)*
- d. *Utilizarea modelelor pentru predicții*

Ca și resursă de bază pentru analiza seriilor de timp utilizăm cartea [4] unde se găsesc elementele de teorie ca și modul de utilizare a sistemului R pentru analiza seriilor de timp. Exemplele sunt în general luate din baze de date din realitate care se instalează odată cu R. Cartea lui Brockwell și Davis¹⁹ conține pas cu pas etapele analizei seriilor de timp cu concentrare pe modelele ARIMA cu explicațiile teoretice necesare și demonstrații pentru teoremele mai ușoare. Exemplele sunt realizate într-un program de calcul elaborat de autori, ITSM2000, care în varianta gratuită este disponibil pe internet (dar este limitat la serii ce nu depășesc 200 de înregistrări). Pentru cei ce urmăresc teoria matematică a seriilor de timp o altă carte a aceluiași autori²⁰, apărută în 1991, pune la dispoziție demonstrații riguroase pentru teoreme dificile din domeniu.

a. Identificarea modelului

Pentru identificarea modelului de serie temporală este nevoie de:

a1. Identificarea tendinței, adică a unei componente ca variază în timp după o formulă cunoscută și care poate fi extrapolată pentru a prezice valori pentru momente viitoare. Uneori unele transformări ale seriei inițiale permit identificarea mai ușoară a tendinței. Identificarea unor componente cu variație periodică exprimabile de asemenea într-o formă care să permită determinarea valorilor lor pentru momente viitoare face ca evoluția deterministă a seriei să fie exprimată prin $m_t + s_t$ unde m_t este tendința neperiodică iar s_t este partea periodică. În anumite determinări ale tendinței partea periodică este inclusă în formula pentru m_t astfel că tendința și partea periodică nu se disting.

a2. După extragerea din serie a tendinței și a componentelor sezoniere urmează analiza resturilor. Dacă aceste resturi se comportă ca niște valori ale unor v. a. independente atunci ele nu mai pot fi prognozate. Dacă însă există o anumită corelație între ele atunci există posibilitatea de a stabili modele stochastice pentru ele

¹⁸ G. E. P. Box, G. M. Jenkins. *Time Series Analysis, Forecasting and Control*. (1976). San Francisco: Holden-Day.

¹⁹ P. J. Brockwell, R. A. Davis. *Introduction to Time Series and Forecasting*. Springer-Verlag New York Berlin Heidelberg, 2002

²⁰ P.J. Brockwell, R.A. Davis, (1991), *Time Series: Theory and Methods*, 2nd Edition, Springer-Verlag, New York.

și a stabili evoluțiile lor cele mai probabile. Pentru aceasta este însă nevoie de considera un model stochastic potrivit.

al. Identificarea tendinței și a componentei periodice

Valorile măsurate ale seriei temporale pot fi privite ca valori ale unui șir de variabile aleatoare $(X_t)_{t \in \mathbb{Z}}$ sau $(X_t)_{t \in \mathbb{N}}$. Descompunerea în tendință plus componenta periodică plus partea reziduală se poate scrie

$$X_t = m_t + s_t + Y_t \quad (7.2.1)$$

În formula de mai sus m_t este tendință iar s_t este componenta periodică. Pentru tendință se utilizează în general metoda celor mai mici pătrate. Se specifică modelul de exemplu $m_t = a + bt$ și se estimează a și b prin metoda celor mai mici pătrate. Pentru specificarea modelului în R avem informații în manualul de introducere în R care vine cu kitul de instalare. Reproducem de acolo câteva specificații (seria de date y este modelată ca expresie de alte serii de date x , x_1 , x_2 , etc., linear, prin formule):

Notăția în R	Explicații
$y \sim x$ sau $y \sim 1 + x$	$y \sim a + bx$
$y \sim 0 + x$ sau $y \sim -1 + x$ sau $y \sim x - 1$	$y \sim bx$
$\log(y) \sim x_1 + x_2$	$\log(y) \sim a + b \cdot x_1 + c \cdot x_2$
$y \sim 1 + x + I(x^2)$	$y \sim a + b \cdot x + c \cdot x^2$
$y \sim 1 + x + I(\sin(2 \cdot x))$	$y \sim a + b \cdot x + c \cdot \sin(2 \cdot x)$

În cadrul seriilor temporale în loc de $y \sim x$ avem tendința ca funcție de timp, $x \sim t$, unde t este timpul.

Determinarea parametrilor a , b , c , etc., se face prin comanda *lm* (linear model) ca mai jos:

```
lm(formula, data, subset, weights, na.action,
    method = "qr", model = TRUE, x = FALSE, y = FALSE, qr =
    TRUE, singular.ok = TRUE, contrasts = NULL, offset, ...)
```

Dintre parametrii funcției *lm* doar “formula” (sub forma unui tabel ca mai sus) este obligatoriu deoarece datele “data” sunt menționate în formulă, iar ceilalți parametri au valori implicite).

Rezultatul funcției *lm* este o listă cu elemente numite din care *coefficients* conține coeficienții obținuți în urma regresiiei liniare, *residuals* conține diferențele între valorile actuale și cele obținute prin formula de regresie. Celelalte componente se pot afla prin *help(lm)* în consola R.

Secvența de mai jos produce un șir de valori care se modelează în două feluri și apoi se reprezintă grafic datele inițiale și modelele:

```
x=1:20/5
y=rnorm(20)/5+2*x+0.2*sin(2*x)
m1=lm(y~1+x)
m2=lm(y~1+x+I(sin(2*x)))
y1=m1$coeff[1]+m1$coeff[2]*x
plot(x,y, main="Modelele y~1+x si y~1+x+I(sin(2*x))")
lines(x,y1, lwd=2, lty=3)
y2=m2$coeff[1]+m2$coeff[2]*x+m2$coeff[3]*sin(2*x)
```

```
lines(x,y2,lwd=3,lty=4)
grid(5,5)
```

Rezultatul este în figura următoare:

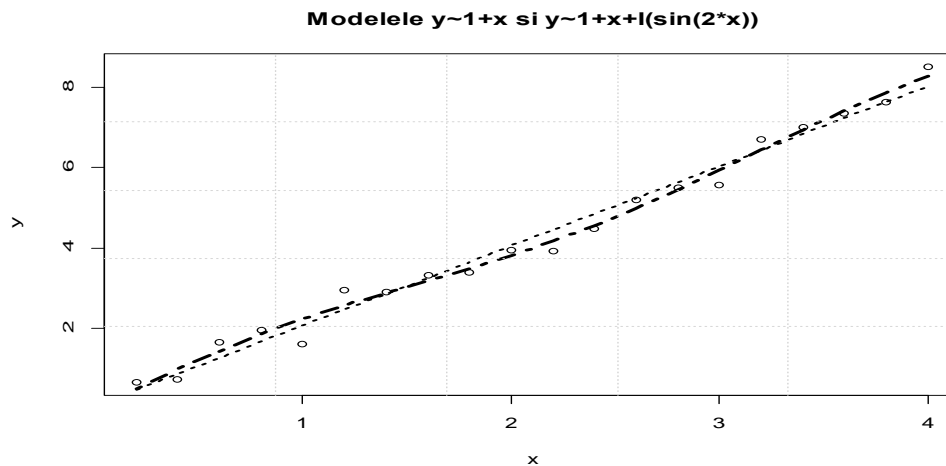


Fig 7.2

Coefficienții celor două modele sunt:

```
> m1

Call:
lm(formula = y ~ 1 + x)

Coefficients:
(Intercept)                x
      0.02623          2.01722

> m2

Call:
lm(formula = y ~ 1 + x + I(sin(2 * x)))

Coefficients:
(Intercept)                x  I(sin(2 * x))
      -0.06758          2.04058          0.26963
```

Pe figură vedem că modelul $y \sim 1+x+I(\sin(2*x))$ aproximează mai bine graficul (x,y) față de modelul $y \sim 1+x$.

Pentru modele neliniare putem utiliza funcția R numită *nlm*. Forma funcției este:

```
nlm(f, p, ..., hessian = FALSE, typsize = rep(1,
length(p)),
      fscale = 1, print.level = 0, ndigit = 12, gradtol =
1e-6,
      stepmax = max(1000 * sqrt(sum((p/typsize)^2)), 1000),
      steptol = 1e-6, iterlim = 100, check.analyticals =
TRUE)
```

Parametrii obligatorii sunt:

- i. f =funcția de minimizat sub forma unei expresii de parametrii căutați și
- ii. p =valori inițiale pentru acești parametri.

De exemplu pentru a modela $y \sim p_1 + p_2 x + p_3 e^{p_4 x}$ unde parametrii p_1, p_2, p_3, p_4 se vor determina prin metoda celor mai mici pătrate trebuie să utilizăm funcția $f(p) = \sum_i (y_i - (p_1 + p_2 x + p_3 e^{p_4 x}))^2$, adică

$$f(p) = \text{sum}((y - p[1] - p[2]*x - p[3]*\exp(p[4]*x))^2)$$

Datele x și y trebuie să fie disponibile la momentul apelării funcției f . Deși nu sunt parametri de intrare pentru f , x și y se văd în contextul în care este funcția f apelată. Scriptul următor face o estimare a parametrilor p .

```
x=1:20/5
y=rnorm(20)/5+2*x+0.2*sin(2*x)
f<-function(p){sum((y-p[1]-p[2]*x-p[3]*exp(p[4]*x))^2)}
m3<-nlm(f,c(1,2,3,0))
p=m3$estimate
y1=p[1]+p[2]*x+p[3]*exp(p[4]*x)
plot(x,y,main="Estimare neliniara")
lines(x,y1)
grid(5,5)
```

Graficul punctelor date și al estimării (regresiei) neliniare este:

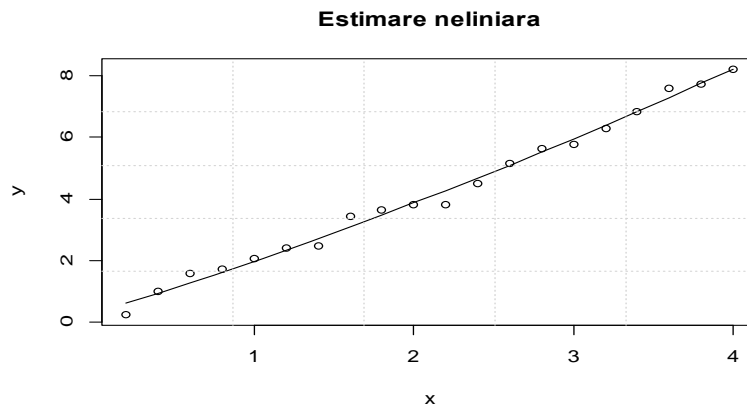


Fig. 7.3

Rezultatul $m3$ al funcției nlm este o listă unde $m3\$estimate$ este un vector cu valorile parametrilor determinați de rutina nlm . Dacă parametrul $\$code$ este 1 sau 2 înseamnă că probabil parametrii p sunt corect estimați (s-a ajuns la concluzia aceasta din motive diferite pentru $\$code=1$ sau $\$code=2$). Vedem mai jos acest rezultat:

```
> m3
$minimum
[1] 2.24936e-14
```

```

$estimate
[1] 4.7903681 6.2172098 0.7584119

$gradient
[1] -2.586661e-06 -2.085010e-06 -4.790367e-05

$code
[1] 2

$iterations
[1] 7

```

Dacă în modelul liniar sau în cel neliniar sunt cuprinse funcții periodice atunci se determină simultan și tendința și partea periodică.

O altă modalitate de identificare a tendinței este sub forma unei *medieri exponențiale*:

Valorile tendinței seriei (x_t) sunt calculate succesiv prin

$$m_t = \alpha x_t + (1 - \alpha)m_{t-1} \quad (7.2.2)$$

cu $\alpha \in (0,1)$. Alegerea lui α este subiectivă. Dacă datele încep cu x_1 atunci avem :

$$m_t = (1 - \alpha)^{t-1} x_1 + \alpha(1 - \alpha)^{t-2} x_2 + \dots + \alpha(1 - \alpha)x_{t-1} + \alpha x_t.$$

În acest fel tendința este păstrată sub forma unui vector de valori $(m_t)_{1 \leq t \leq T}$.

Estimarea după timpul maxim de măsurare, T , se face prin $m_{T+k} = m_T$ ceea ce nu este foarte convenabil. O estimare a tendinței în așa fel ca ea să poată fi extinsă într-un mod mai rezonabil după $t=T$ se poate face prin procedura Holt care constă în mediarea exponențială a tendinței m_t cât și a pantei P_t sub forma:

$$\begin{aligned} m_1 &= x_1 \\ P_1 &= \frac{(x_n - x_1)}{n-1} \text{ pentru un } n, \text{ de exemplu } n=2 \\ m_t &= \alpha x_t + (1 - \alpha)(m_{t-1} + P_{t-1}) \\ P_t &= \beta(m_t - m_{t-1}) + (1 - \beta)P_{t-1} \end{aligned} \quad (7.2.3)$$

Parametrii α, β sunt subiectiv aleși în intervalul $(0,1)$. Extrapolarea tendinței pentru $t > T$ se face linear prin

$$m_{T+k} = m_T + P_T \cdot k \quad (7.2.4)$$

O procedură mai generală este Holt-Winters aditiv care mediază exponențial tendința m_t , componenta sezonieră s_t cât și panta tendinței P_t . Pentru aceasta este nevoie de trei parametri α, β, γ aleși subiectiv în intervalul $(0,1)$. De asemenea este nevoie apriori de lungimea L a părții periodice și de valorile inițiale pentru partea periodică : $(s_t)_{t=1..L}$.

$$\begin{aligned} m_t &= \alpha(x_t - s_{t-L}) + (1 - \alpha)(m_{t-1} + P_{t-1}) \\ P_t &= \beta(m_t - m_{t-1}) + (1 - \beta)P_{t-1} \\ s_t &= \gamma(x_t - m_t) + (1 - \gamma)s_{t-L} \end{aligned} \quad (7.2.5)$$

Predicția valorilor dincolo de pragul $t=T$ se face după formula:

$$x_{T+k} = m_T + k * P_T + s_{T+1+(h-1) \bmod L} \quad (7.2.6)$$

Există și un prececeu Holt-Winters multiplicativ.

Procedura Holt-Winters implementată în R este apelabilă prin:

```
HoltWinters(x, alpha = NULL, beta = NULL, gamma = NULL,
seasonal = c("additive", "multiplicative"), start.periods
= 2, l.start = NULL, b.start = NULL, s.start = NULL,
optim.start = c(alpha = 0.3, beta = 0.1, gamma =
0.1), optim.control = list())
```

Dintre parametrii funcției doar x de tip *ts* este obligatoriu. Rezultatul funcției este un obiect *HoltWinters* adică o listă cu componentele *fitted*, *x*, *alpha*, *beta*, *gamma*, *coefficients*, *seasonal*, *SSE*, *call*. Lungimea L a perioadei este luată egală cu frecvența (frequency) utilizată la definirea seriei temporale. Dintre elementele listei rezultat, *fitted* este o listă ce conține timpul t , valoarea estimată \hat{x}_t , m_t numit level, P_t numit tendință și s_t . Prin comanda *help(HoltWinters)* se poate afla și semnificația celorlalți parametri.

Predicția valorilor ce urmează în serie după timpul maxim T de măsurare se poate face prin funcția *predict*. În scriptul ce urmează după determinarea tendinței și a părții sezoniere (perioada=168) pentru datele CO dinainte, se face o predicție pentru o săptămână (=168 ore) pentru nivelul de CO. Mai jos avem graficul concentrației de CO cu negru (seria de timp *tsco* a fost creată mai înainte) împreună cu valorile netezite prin metoda Holt-Winters cu verde și predicția Holt-Winters pentru 168 de ore cu albastru.

```
plot(tsco, main="Co masurat + CO prezis Holt-
Winters", cex=2)
m4<- HoltWinters(tsco)
lines(fitted(m4)[,1], col = 3)
pr <- predict(m4, 168, prediction.interval = FALSE)
lines(pr, main="Predictie Holt-Winters", col=4)
```

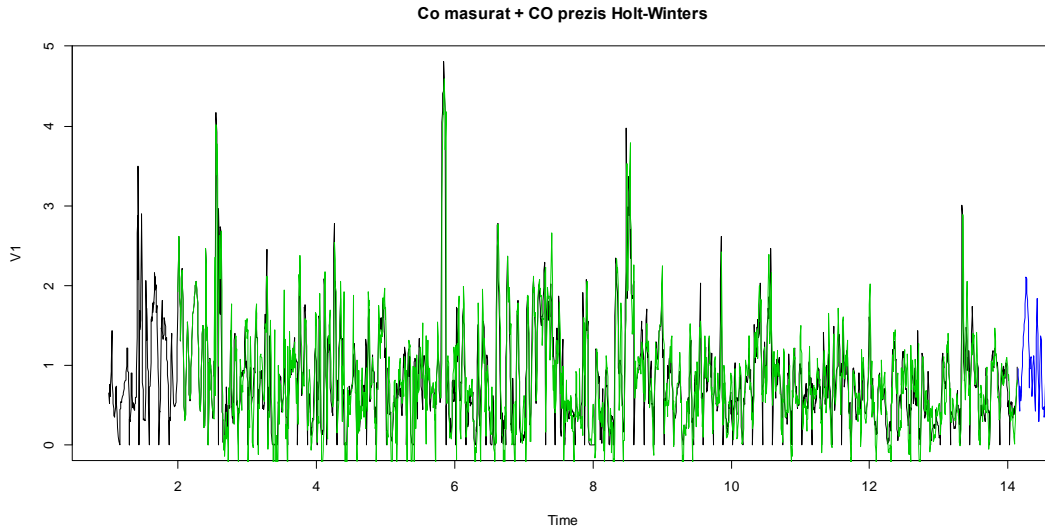


Fig. 7.4

Funcția *predict* are forma:

```
predict(object, n.ahead=1, prediction.interval = FALSE,
        level = 0.95, ...)
```

unde *object* este un model (de exemplu Holt-Winters) ce conține elementele de predicție după ultima valoare măsurată (conform formulelor Holt-Winters de mai sus). Rezultatul funcției *predict* este o serie de timp cu valorile prezise în continuarea celor ale seriei inițiale.

O altă modalitate de eliminare a tendinței este procedura de diferențiere. Astfel de la o serie (x_t) se obține o nouă serie (y_t) cu $y_t = x_{t+1} - x_t$. O nouă diferențiere conduce la seria (z_t) cu $z_t = y_{t+1} - y_t = x_{t+2} - 2x_{t+1} + x_t$ care se numește diferența de ordin doi s.a.m.d. Funcția *diff* din R realizează aceste diferențe finite la dreapta:

```
diff(x, lag = 1, differences = 1, ...)
```

Parametrii sunt x =vector sau matrice de numere, *differences* este ordinul diferenței finite, *lag* este pasul de timp (de exemplu la *lag*=2 prima diferență se calculează după formula $y_t = x_{t+2} - x_t$). Dacă notăm B operatorul de întârziere $(BX)_t = X_{t-1}$ atunci seria obținută prin prima diferență la *lag*=1 este $Y = (1 - B)X$, iar la *lag*=2 este $Y = (1 - B^2)X$. Dacă $y = \text{diff}(x)$ atunci x poate fi restaurat prin $c(x[1], x[1] + \text{cumsum}(y))$.

Transformarea seriei într-o serie staționară (vezi mai jos) se poate încerca și printr-o transformare de tip Box-Cox:

$$y_t = \begin{cases} \frac{x_t^\lambda - 1}{\lambda} & \text{pour } \lambda \neq 0 \\ \ln(x_t) & \text{pour } \lambda = 0 \end{cases} \quad (7.2.7)$$

În pachetul *TeachingDemos* există funcția *vis.bboxcox* care permite vizualizarea interactivă a graficului seriei temporale după o transformare Box-Cox. De asemenea funcția *bct(y, lambda)* din același pachet calculează transformarea Box-Cox de mai sus.

a2. Determinarea modelului pentru reziduali

Analiza rezidualilor (diferențele dintre valorile măsurate și cele date de tendință și componenta sezonieră) are scopul de a determina dacă partea care rămâne după extragerea tendinței și a periodicității este formată din valori fără legătură între ele, caz în care aceste valori reziduale nu mai pot fi prezise, sau au o anumită corelație și atunci se poate încerca o previziune a lor bazată pe un model statistic. Se pot încerca testele *Box.test* (din pachetul *stats*) sau *runs.test* (din pachetul *tseries*). În cele ce urmează ne vom ocupa de serii de timp staționare. Un șir de variabile aleatoare $(X_t)_{t \in \mathbb{Z}}$ se numește temporală staționară slab (staționar de ordinul 2) dacă media $\mu_X(t) = E(X_t)$ este constantă (nu depinde de t) și funcția de *autocovarianță*

$$\gamma_X(r, s) = E[(X_r - \mu_X(r))(X_s - \mu_X(s))] \quad (7.2.8)$$

depinde doar de diferența $h=r-s$. În cele ce urmează notăm $E(X)$ media v.a. X . Pentru o serie temporală staționară notăm $\gamma_X(h) = E((X_{t+h} - \mu_X)(X_t - \mu_X))$, independent de t . Avem:

$$\begin{aligned} \gamma_X(0) &= \text{var}(X_t) > 0 \\ \gamma_X(h) &= \gamma_X(-h) \\ |\gamma_X(h)| &\leq \gamma_X(0) \\ \forall k \geq 1, \forall a_1, a_2, \dots, a_k &\Rightarrow \sum_{i,j=1}^k a_i \gamma(i-j) a_j \geq 0 \end{aligned}$$

În general vom scrie $\gamma(h)$ sau γ_h în loc de $\gamma_X(h)$ dacă nu este pericol de confuzie asupra seriei temporale X . Numim funcția de *autocorelație* funcția:

$$\rho(h) = \frac{\gamma(h)}{\gamma(0)} \quad (7.2.9)$$

În general vom scrie ρ_h în loc de $\rho(h)$ Estimatori numerici pentru medie, funcția de autocovarianță și autocorelație sunt:

$$\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n} \quad (7.2.10)$$

$$\hat{\gamma}_h = \frac{1}{n} \sum_{t=1}^{n-|h|} (X_{t+|h|} - \bar{X}_n)(X_t - \bar{X}_n) \quad (7.2.11)$$

$$\hat{\rho}_h = \frac{\hat{\gamma}_h}{\hat{\gamma}_0} \quad (7.2.12)$$

Box și Jenkins recomandă $n \geq 50$ et $h \leq n/4$ pentru estimările de mai sus.

Estimatorii pentru medie și corelație pot fi studiați doar în unele ipoteze asupra seriei temporale X . Astfel introducem unele tipuri de serii temporale staționare :

Zgomotul alb este un tip de serie temporală $(Z_t)_{t \in \mathbb{Z}}$ în care fiecare v.a. este de medie zero și varianță σ^2 , iar covarianța este

$$\gamma_Z(h) = \begin{cases} \sigma^2, & h = 0 \\ 0, & h \neq 0 \end{cases} \quad (7.2.13)$$

Variabilele aleatoare Z_t nu sunt neapărat independente.

Seria temporală i.i.t. este la fel ca zgomotul alb dar toate variabilele Z_t sunt independente identic distribuite.

Modelul general linear este o serie temporală $(X_t)_{t \in \mathbb{Z}}$ în care $X_t = \sum_{j=-\infty}^{\infty} \psi_j Z_{t-j}$,

$(Z_t)_{t \in \mathbb{Z}}$ este un zgomot alb și $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$. Se verifică faptul că în acest caz avem

$$\gamma_X(h) = \sigma^2 \sum_{j=-\infty}^{\infty} \psi_{j+h} \psi_j.$$

Serii staționare MA(q) (medii mobile) sunt serii de forma $X_t = Z_t + \sum_{i=1}^q \theta_i Z_{t-i}$

unde $(Z_t)_{t \in \mathbb{Z}}$ este un zgomot alb. Notăm $\theta_0 = 1$. Avem $E(X_t) = 0$ și

$$\gamma_X(h) = \text{Cov}(X_{t+h}, X_t) = \begin{cases} \sigma^2 \sum_{i=0}^{q-|h|} \theta_i \theta_{i+|h|}, & |h| \leq q \\ 0, & |h| > q \end{cases} \quad (7.2.14)$$

Seriile staționare de tip MA(q) sunt singurele serii staționare care au $\gamma_X(h) = 0$ pentru $h > q$.

Seria staționară autoregersivă AR(p) este o serie de tipul (în general cerem și ca media să fie zero) $X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + Z_t$ unde $(Z_t)_{t \in \mathbb{Z}}$ este un zgomot alb și Z_t este independent de X_{t-k} , $k > 0$. Dacă polinomul

$\phi(z) = 1 - \phi_1 z - \phi_2 z^2 - \dots - \phi_p z^p$ are toate rădăcinile mai mari ca 1 în modul (serie cauzală) atunci $X_t = Z_t + \sum_{j=1}^{\infty} \psi_j Z_{t-j}$ cu $\sum_{j=1}^{\infty} |\psi_j| < \infty$. Inmulțind relația de definiție a lui X_t cu X_{t-k} găsim sistemul și luând mediile obținem pentru $k=1, 2, \dots, p$:

$$\begin{cases} \rho_1 = \phi_1 + \phi_2 \rho_1 + \phi_3 \rho_2 + \dots + \phi_p \rho_{p-1} \\ \rho_2 = \phi_1 \rho_1 + \phi_2 + \phi_3 \rho_1 + \dots + \phi_p \rho_{p-2} \\ \dots \\ \rho_p = \phi_1 \rho_{p-1} + \phi_2 \rho_{p-2} + \phi_3 \rho_{p-3} + \dots + \phi_p \end{cases} \quad (7.2.15)$$

care prin rezolvare ne dă $\rho_1, \rho_2, \dots, \rho_p$. Avem de asemenea $\rho_0 = 1$. Pentru $k > p$ avem analog $\rho_k = \phi_1 \rho_{k-1} + \phi_2 \rho_{k-2} + \phi_3 \rho_{k-3} + \dots + \phi_p \rho_{k-p}$. Funcția de autocorelație se poate astfel determina complet și nu se anulează în principiu nicăieri.

Se definește funcția de autocorelație parțială $\phi_{k,k}$ în felul următor: Pentru X_t și X_{t+k} se consideră componentele perpendiculare pe spațiul vectorial generat de $X_{t+1}, X_{t+2}, \dots, X_{t+k-1}$, fie ele $X_{t,0}$ respectiv $X_{t,k}$. Funcția de autocorelație parțială se definește acum ca și $\phi_{k,k} = \text{Corr}(X_{t,0}, X_{t,k})$. Deoarece proiecțiile depind doar de produsele scalare reciproce (corelații între X_i, X_j) și seria este staționară rezultă că $\phi_{k,k}$ nu depinde de t ci doar de ρ_i pentru $i=0,1,2,\dots,k$. Pentru seriile autoregresive $\phi_{k,k}$ este zero dacă $k > p$ și aceasta este o modalitate de a detecta dacă o serie este de tip AR(p) estimând pe $\phi_{k,k}$ prin $\hat{\phi}_{k,k}$ ce se calculează cu ajutorul estimatorilor $\hat{\rho}_i$ pentru corelațiile ρ_i .

Serii staționare autoregresive medie mobilă ARMA(p,q) sunt serii care verifică o relație de forma:

$$X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \dots + \theta_q Z_{t-q} + Z_t$$

unde Z_t este un zgomot alb independent de X_{t-k} pentru $k > 0$. Dacă seria este cauzală adică polinomul $\phi(z) = 1 - \phi_1 z - \phi_2 z^2 - \dots - \phi_p z^p$ are rădăcinile mai mari ca 1 în modul

atunci $X_t = Z_t + \sum_{j=1}^{\infty} \psi_j Z_{t-j}$ cu $\sum_{j=1}^{\infty} |\psi_j| < \infty$. Dacă și polinomul

$\theta(z) = 1 + \theta_1 z + \theta_2 z^2 + \dots + \theta_q z^q$ are toate rădăcinile în modul mai mari ca 1 seria se numește invertibilă. Funcția de autocovarianță se poate determina exact pentru seriile cauzale dacă se dau coeficienții ϕ, θ, σ^2 , unde σ^2 este varianța lui Z_t .

Funcția ARMAacf din pachetul stats face acest lucru în R. Apelul se face prin:

```
ARMAacf(ar = numeric(0), ma = numeric(0), lag.max = r, pacf = FALSE)
```

Se vede din apelul funcției că prin opțiunea *pacf* putem obține funcția de autocorelație parțială.

Serii nestacionare ARIMA(p,d,q) sunt seriile care prin diferența de ordin d la dreapta devine staționară de tip ARMA(p,q).

Am văzut în modelele de serii temporale staționare că tipurile MA(q), AR(p) sunt determinate de corelațiile ρ_k . Estimatorii pentru Proprietățile estimatorilor pentru medie și corelație pentru seriile staționare pot fi enumerate astfel:

i. Dacă $\sum_{h=-\infty}^{\infty} |\gamma_X(h)| < \infty$ și $v = \sum_{h=-\infty}^{\infty} \gamma_X(h)$ atunci $\bar{X}_n - \mu$ este $AN(0, \frac{v}{n})$ (AN=asimptotic normală). În particular acest lucru se întâmplă dacă $X_t = \mu + \sum_{j=-\infty}^{\infty} \psi_j Z_{t-j}$ cu $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$ și $(Z_t)_{t \in \mathbb{Z}}$ este un zgomot alb.

ii. Dacă $X_t = \mu + \sum_{j=-\infty}^{\infty} \psi_j Z_{t-j}$ cu $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$, $(Z_t)_{t \in \mathbb{Z}}$ este un zgomot alb și în plus

$$\sum_{j=-\infty}^{\infty} |j| \psi_j^2 < \infty, \text{ atunci } \sum_{h=-\infty}^{\infty} |\gamma_X(h)| = \sigma^2 \left(\sum_{j=-\infty}^{\infty} |\psi_j| \right)^2 < \infty \text{ și}$$

$$\begin{pmatrix} \hat{\rho}(1) \\ \hat{\rho}(2) \\ . \\ . \\ \hat{\rho}(m) \end{pmatrix} \text{ est AN } \begin{pmatrix} \rho(1) \\ \rho(2) \\ . \\ . \\ \rho(m) \end{pmatrix}, \frac{1}{n} W$$

unde matricea $W \in \mathbb{R}^{m \times m}$ are componentele

$$\begin{aligned} w_{i,j} &= \sum_{k=-\infty}^{\infty} \left[\rho(k+i)\rho(k+j) + \rho(k-i)\rho(k+j) + 2\rho(i)\rho(j)\rho^2(k) - \right. \\ &\quad \left. - 2\rho(i)\rho(k)\rho(k+j) - 2\rho(j)\rho(k)\rho(k+i) \right] = \\ &= \sum_{k=1}^{\infty} (\rho(k+i) + \rho(k-i) - 2\rho(i)\rho(k)) \cdot (\rho(k+j) + \rho(k-j) - 2\rho(j)\rho(k)) \end{aligned} \quad (7.2.16)$$

(formula lui Bartlett).

În particular $n^{1/2}(\hat{\rho}(h) - \rho(h))$ are pentru n mare o distribuție normală de medie 0 și varianță $w_{h,h}$.

În R funcția de autocorelație empirică $\hat{\rho}$ este calculată cu funcția *acf* care are forma:

```
acf(x, lag.max = NULL,
    type = c("correlation", "covariance", "partial"),
    plot = TRUE, na.action = na.fail, demean = TRUE, ...)
```

Rezultatul aplicării funcției *acf* este o listă ce conține *lag*=intervale unde este estimat $\hat{\rho}$, *acf*=estimările pentru $\hat{\rho}$, *type*=la fel ca în apelul funcției *acf*, *n.used*=numărul de observații din serie, *series*=numele seriei de timp, *snames*=numele componentelor în o serie multivariată. Se vede din definiție că funcția *acf* determină și corelația parțială dacă alegem corespunzător parametrul *type*. Dacă alegem *plot=TRUE* atunci este tipărită funcția de autocorelație împreună cu intervalul de încredere de 95%.

Pentru determinarea tipului de serie temporală între tipurile MA(q), AR(p), ARMA(p,q) procedăm astfel:

MA. Dacă funcția de autocorelație dată de *acf* este nulă pentru intervale temporale de lungime $k > q$ (sau cade în intervalul de încredere cu încrederea 95% în jurul lui zero, tipărit pe grafic de funcția *acf*), atunci alegem modelul MA(q).

AR. Dacă funcția de autocorelație parțială dată de *acf* este nulă pentru intervale temporale $k > p$ (sau cade în intervalul de încredere cu încrederea 95% în jurul lui zero, tipărit pe grafic de funcția *acf*), atunci alegem modelul AR(*p*).

Putem simula în R serii temporale de tip ARIMA(*p,d,q*) prin funcția *arima.sim* care are forma:

```
arima.sim(model, n, rand.gen = rnorm, innov = rand.gen(n, ...),
           n.start = NA, start.innov = rand.gen(n.start, ...),...)
```

Doar parametrii model și n sunt obligatorii, ceilalți având valori implicite. Mai jos simulăm o serie AR(2) și îi tipărim ACF și PACF. Scriptul este:

```
s1<-arima.sim(model=list(ar=c(-.3,-.4)), n=10000)
par(mfrow=c(1,2))
acf(s1)
acf(s1, type="partial")
```

Rezultatul este în figura următoare:

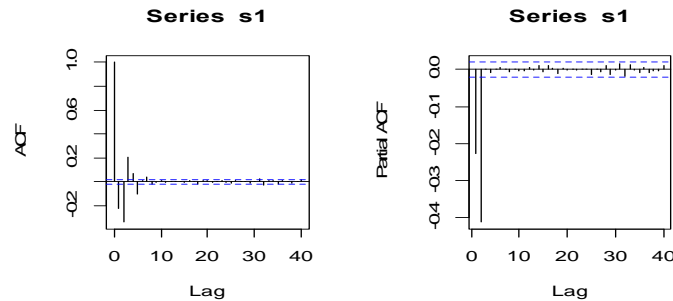


Fig. 7.5

Forma PACF sugerează un model AR(2). Acum simulăm un model MA(2) și tipărim ACF+PACF:

```
s1<-arima.sim(model=list(ma=c(-.3,-.4)), n=10000)
par(mfrow=c(1,2))
acf(s1)
acf(s1, type="partial")
```

Rezultatul este:

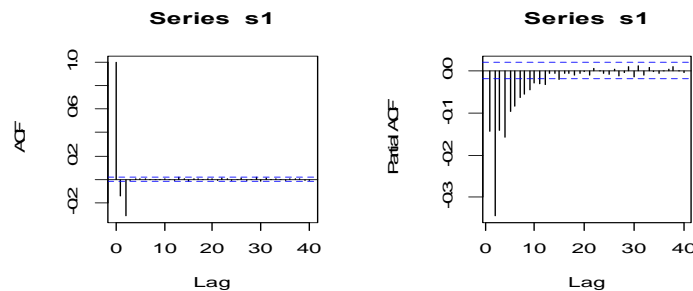


Fig. 7.6

Graficele sugerează un model MA(2).

ARMA. Pentru modelele ARMA(p,q) un indiciu asupra valorilor potrivite pentru p și q este funcția extinsă de autocorelație²¹ (vezi [4], pag. 116). Funcția extinsă de autocorelație EACF se calculează pentru parametrii x =serie temporală, p,q=numere naturale în intervalele [0, pmax] respectiv [0, qmax]. Rezultatul funcției este o mulțime de coeficienți calculați după un algoritm complicat dar care ar trebui să formeze în octantul unde p și q sunt coordonate un triunghi de zerouri cu vârful în valorile actuale pentru p și q ai seriei ARMA(p,q). Funcția *eacf* se găsește în pachetul TSA și are forma:

```
eacf(z, ar.max = 7, ma.max = 13)
```

Apelată pe o serie ARMA(2,2) simulată ne dă următorul rezultat:

```
s1<-arima.sim(model=list(ma=c(-.3,.5), ar=c(0.4,-0.2)),
n=10000)
par(mfrow=c(1,2))
acf(s1)
acf(s1, type="partial")
eacf(s1)
```

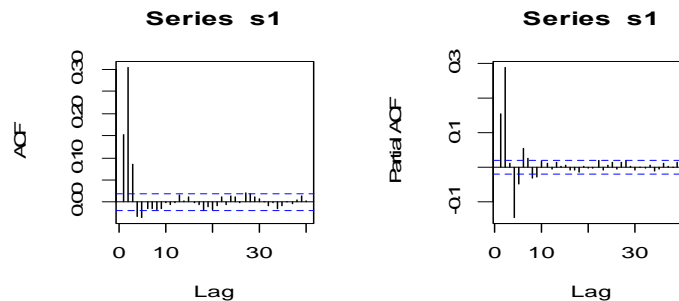


Fig. 7.7

```
> eacf(s1)
AR/MA
 0 1 2 3 4 5 6 7 8 9 10 11 12 13
0 x x x x x 0 0 0 0 0 0 0 0 0
1 x x x x 0 0 0 0 0 0 0 0 0 0
2 x x 0 x x 0 0 0 0 0 0 0 0 0
3 x x x 0 0 0 0 0 0 0 0 0 0 0
4 x x x x 0 0 x 0 0 0 0 0 0 0
5 x x x 0 x 0 x 0 0 0 0 0 0 0
6 x x x 0 x x 0 x 0 0 0 0 0 0
7 x x 0 x x x x x 0 x 0 0 0 0
```

Pe graficele ACF și PACF seria nu pare AR sau MA. EACF sugerează un model ARMA(2,2), ARMA(1,4) sau MA(6). Dacă se rulează exemplul din nou se obține altă figură dată de *eacf* din cauză că simularea produce alte valori pentru serie.

O altă metodă de determinare a parametrilor p și q este de a propune toate valorile posibile pentru p și q în intervalele [0, pmax] respectiv [0, qmax] ca în cazul

²¹ Tsay, R. S. and Tiao, G. (1984). "Consistent estimates of autoregressive parameters and extended sample autocorrelation function for stationary and nonstationary ARMA Models." Journal of the American Statistical Association, 79, 385, 84–96.

ea \acute{c} f și de a determina coeficienții modelului (vezi mai departe). Valorile p și q optime se determină prin criteriile AIC, BIC etc. (vezi mai departe).

b. Estimarea parametrilor

După determinarea tipului ARMA(p,q) urmează estimarea parametrilor. Pentru modelele AR(p) *metoda momentelor* (Yule-Walker) este eficace. Această metodă constă în rezolvarea sistemului (7.2.15) în necunoscutele $\phi_1, \phi_2, \dots, \phi_p$ dacă $\rho_1, \rho_2, \dots, \rho_p$ se înlocuiesc cu estimările lor $\hat{\rho}_1, \hat{\rho}_2, \dots, \hat{\rho}_p$. Pentru modele ce conțin și parte medie mobilă metoda este mult mai complicată.

Metoda celor mai mici pătrate constă în a scrie:

$$Z_t = X_t - \phi_1 X_{t-1} - \phi_2 X_{t-2} - \dots - \phi_p X_{t-p} - \theta_1 Z_{t-1} - \theta_2 Z_{t-2} - \dots - \theta_q Z_{t-q} \quad (7.2.17)$$

și a minimiza:

$$S(\phi, \theta) = \sum_{t=p+1}^T (Z_t)^2 \quad (7.2.18)$$

considerând $Z_p, Z_{p-1}, \dots, Z_{p-q+1}$ egale cu 0. Metoda este aplicabilă și pentru serii multidimensionale.

Metoda cea mai utilizată este *metoda verosimilității maxime* deși metoda celor mai mici pătrate a lui Hannan și Rissanen poate fi îmbunătățită încât să aibă aceeași eficiență asimptotică. Funcția de verosimilitate are o expresie complicată și poate fi calculată cu ajutorul algoritmului inovațiilor. Detalii se pot găsi în cărțile menționate ale lui Brockwell și Davis de exemplu²², demonstrațiile găsindu-se în cartea pentru matematicieni²³. Tot prin metoda verosimilității maxime se poate stabili care sunt valorile optime pentru parametrii p și q în modelul ARMA. Pentru aceasta se estimează parametrii $\phi_1, \phi_2, \dots, \phi_p, \theta_1, \dots, \theta_q$ și σ^2 prin metoda verosimilității maxime și apoi se calculează criteriul AICC:

$$AICC = -2 \log(L(\phi_1, \phi_2, \dots, \phi_p, \theta_1, \dots, \theta_q, \sigma^2)) + 2(p + q + 1)/(n - p - q - 2) \quad (7.2.19)$$

În acest criteriu $L(\phi_1, \phi_2, \dots, \phi_p, \theta_1, \dots, \theta_q, \sigma^2)$ este funcția de verosimilitate a cărei expresie

$$L(\phi, \theta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} (r_0 r_1 \dots r_{T-1}) e^{-\frac{1}{2\sigma^2} \sum_{j=1}^T \frac{(x_j - \hat{x}_j)^2}{r_{j-1}}} \quad (7.2.20)$$

depinde $\hat{X}_1, \hat{X}_2, \dots, \hat{X}_T$ și r_0, r_1, r_{T-1} calculați prin algoritmul inovațiilor (vezi citările de la subsol). Criteriul constă în a alege parametrii p și q care minimizează AICC. Penalizarea $2(p + q + 1)/(n - p - q - 2)$ determină ca valorile optime pentru p și q să nu fie foarte mari.

În sistemul R există funcții care determină parametrii pentru modelele AR sau ARIMA (deci și ARMA). Pentru modelul AR avem funcția:

²² P. J. Brockwell, R. A. Davis, *Introduction to Time Series and Forecasting*. Springer-Verlag New York Berlin Heidelberg, 2002, pag. 156-162

²³ P.J. Brockwell, R.A. Davis, (1991), *Time Series: Theory and Methods*, 2nd Edition, Springer-Verlag, New York, cap. 8

```
ar(x, aic = TRUE, order.max = NULL,
    method=c("yule-walker", "burg", "ols", "mle", "yw"),
    na.action, series, ...)
```

Metoda "yule-walker" este metoda momentelor, "ols" este metoda celor mai mici pătrate, iar "mle" este metoda verosimilității maxime. Sunt calculate modelele AR până la ordinul order.max și se alege modelul cu aic minim. Seria de timp x poate fi și multidimensională.

Rezultatul funcției *ar* este o listă ce conține: *order*, *ar*, *var.pred*, *x.mean*, *x.intercept*, *aic*, *n.used*, *partialacf*, *order.max*, *resid*, *method*, *series*, *call*, *asy.var.coef*. Cu comanda *help(ar)* se pot obține informații asupra semnificației acestor elemente. Avem de exemplu *order*=ordinul p al modelului AR determinat, *resid*=diferențele $X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p}$, *ar*=coeficienții ϕ_1, \dots, ϕ_p estimați.

Iată un model de aplicare a funcției *ar*. Se simulează un model AR și apoi se estimează parametrii săi fără a cunoaște lungimea p a lor. Scriptul este mai jos:

```
s2<-arima.sim(model=list(ar=c(0.4,-0.2)), n=1000)
m1<-ar(s2,order.max=5,AIC=T,method='yw')
m2<-ar(s2,order.max=5,AIC=T,method='ols')
m3<-ar(s2,order.max=5,AIC=T,method='mle')
```

Rezultatele se văd mai jos în consola R:

```
> m1
```

```
Call:
```

```
ar(x = s2, order.max = 5, method = "yw", AIC = T)
```

```
Coefficients:
```

```
      1      2
0.390  -0.209
```

```
Order selected 2  sigma^2 estimated as  1.029
```

```
> m2
```

```
Call:
```

```
ar(x = s2, order.max = 5, method = "ols", AIC = T)
```

```
Coefficients:
```

```
      1      2
0.3905  -0.2092
```

```
Intercept: -0.0003215 (0.03208)
```

```
Order selected 2  sigma^2 estimated as  1.027
```

```
> m3
```

```
Call:
```

```
ar(x = s2, order.max = 5, method = "mle", AIC = T)
```

```
Coefficients:
```

```
      1      2
```



```
0.3901 -0.2089
```

```
Order selected 2 sigma^2 estimated as 1.026
```

Vedem că a fost determinat corect ordinul $p=2$ și s-au determinat aproximativ corect coeficienții.

Pentru estimarea modelelor ARMA(p,q) sau ARIMA(p,d,q) avem funcția *arima* în pachetul stats, care are forma:

```
arima(x, order = c(0, 0, 0),
      seasonal = list(order = c(0, 0, 0), period = NA),
      xreg = NULL, include.mean = TRUE,
      transform.pars = TRUE,
      fixed = NULL, init = NULL,
      method = c("CSS-ML", "ML", "CSS"),
      n.cond, optim.method = "BFGS",
      optim.control = list(), kappa = 1e6)
```

Parametrii sunt x : o serie temporală univariată sau un vector numeric, *order*: specificare a tipului de modelare, *seasonal*: specificare a părții sezoniere, *xreg*: variabile externe față de care se face o regresie liniară și se modelează apoi ARIMA rezidualii, *method*: metoda de determinare a parametrilor (implicită "CSS"=metoda celor mai mici pătrate), *optim*: este metoda utilizată pentru minimizare (implicită="BFGS", adică Broyden, Fletcher, Goldfarb and Shanno). Mai multe detalii cu `help(arima)`. Rezultatul funcției ARIMA este o listă *arima* ce conține:

- coef*: o listă cu coeficienții
- sigma2*: varianța σ^2 a variabilelor Z_t din modelul ARMA(p,q)
- residuals*: valorile Z_t așa cum rezultă din estimare
- arma*: specificarea modelului
- aic*: coeficientul Akaike, valabil doar pentru metoda ML (ver. maximă)
- code*: codul de convergență întors de funcția *optim*; dacă *code*=0 atunci optimizarea a decurs normal, iar pentru alte coduri se poate utiliza comanda `help(optim)`
- series*: numele seriei x

Ceilalți parametri de ieșire pot fi consultați cu `help(arima)`.

Următorul script generează o seria ARMA(2,2) și apoi se încearcă cu funcția *arima* să se determine modelul.

```
s1<-arima.sim(model=list(ma=c(-.3,.5), ar=c(0.4,-0.2)),
n=10000)
modell<-arima(s1,order=c(2,0,2),method="ML")
```

Rezultatul este:

```
> modell
```

Call:

```

arima(x = s1, order = c(2, 0, 2), method = "ML")

Coefficients:
          ar1          ar2          ma1          ma2  intercept
    0.4012   -0.2091   -0.3021    0.4996   -0.0040
s.e.    0.0310    0.0295    0.0279    0.0243    0.0147

sigma^2 estimated as 0.9877:  log likelihood = -14127.79,
aic = 28265.57

```

Se vede că practic au fost determinați corect coeficienții ar și ma . Coeficientul *intercept* din lista de coeficienți este media variabilelor X_t ale seriei de timp (probabil ar trebui să se schimbe denumirea în *mean*).

O extensie a funcției *arima* din pachetul stats se găsește în pachetul TSA cu numele tot *arima*. Cititorul interesat poate găsi informațiile cu ajutorul comenzii *help* din R.

Până la versiunea 2.14.0 nu există variante de modelare pentru serii de timp ARIMA(p,d,q) multidimensionale. De asemenea trebuie menționat că modelul ARMA trebuie să fie cauzal (modelele necauzale sunt echivalente cu modele cauzale dar cu alt zgomot alb). De asemenea se poate considera modelul invertibil pentru că pentru orice model neinvertibil există unul invertibil cu aceeași funcție de autocorelație (care determină în final coeficienții modelului). Detalii se găsesc în cărțile menționate ale lui Brockwell și Davis.

c. Validarea modelului (diagnosticul veridicității modelului ales)

Pentru validarea modelului trebuie verificat că rezidualii produși de modelele statistice sunt variabile aleatoare independente de medie zero și aceeași varianță σ^2 . De regulă se fac teste pentru a verifica dacă valorile reziduale Z_t provin din v.a. independent identic distribuite. Pentru aceasta există un număr de teste.

1. *Funcția de autocorelație* $\hat{\rho}_n$ ar trebui să se anuleze pentru $n \geq 1$. În R acest lucru este cercetat prin funcția *acf*.

2. *Testul rangului*. Dacă y_1, y_2, \dots, y_n sunt realizări ale unei serii i.i.d. atunci fie P numărul de perechi (i, j) astfel ca $i > j$ și $y_i > y_j$. Pentru n mare P are o repartiție aproximativ $N(m, \sigma^2)$ cu $m = \frac{n(n-1)}{4}$, $\sigma^2 = n(n-1)(2n+5)/72$

3. *Testul Ljung-Box* constă în a calcula: $Q(h) = n(n+2) \sum_{t=1}^h \frac{\hat{\rho}(t)}{n-t}$. Pentru n mare și y_i v.a. iid, $Q(h)$ are o repartiție aproximativ χ^2 cu h grade de libertate.

4. *Testul semnelor diferențelor*. În acest test S este numărul de valori i pentru care $y_i > y_{i-1}$. Pentru n mare și valorile y_i provenind de la v.a. iid, S are o repartiție aproximativ $N(m, \sigma^2)$ cu $m = \frac{n-1}{2}$, $\sigma^2 = \frac{n+1}{12}$.

5. *Testul punctelor de întoarcere.* La momentul i avem un punct de întoarcere dacă $y_{i-1} < y_i > y_{i+1}$ sau $y_{i-1} > y_i < y_{i+1}$. Fie T numărul punctelor de întoarcere. Dacă n este mare și valorile y_i provin de la v.a. iid atunci T are o repartiție aproximativ $N\left(\frac{2(n-2)}{3}, \frac{16n-29}{90}\right)$

În R există funcția *tsdiag* care tipărește rezidualii standardizați (la varianță 1), calculează funcția de autocorelație a rezidualilor și calculează testul Box (sau Ljung–Box) până la un h dat. Forma funcției este:

```
tsdiag(object, gof.lag, ...)
```

unde *object* este o serie temporală modelată în R (deci are rezidualii calculați), *gof.lag* este numărul maxim de pași h (vezi testul 3 (Ljung–Box) până la care se calculează $Q(h)$).

În scriptul următor se simulează o serie ARMA(2,2) căreia i se determină parametrii cu funcția *arima* și apoi cu funcția *tsdiag* se studiază dacă modelul este valid.

```
s1<-arima.sim(model=list(ma=c(-.3,.5),      ar=c(0.4,-0.2)),
n=10000)
modell<-arima(s1,order=c(2,0,2),method="ML")
tsdiag(modell)
```

Obținem următoarele reprezentări grafice din *tsdiag*:

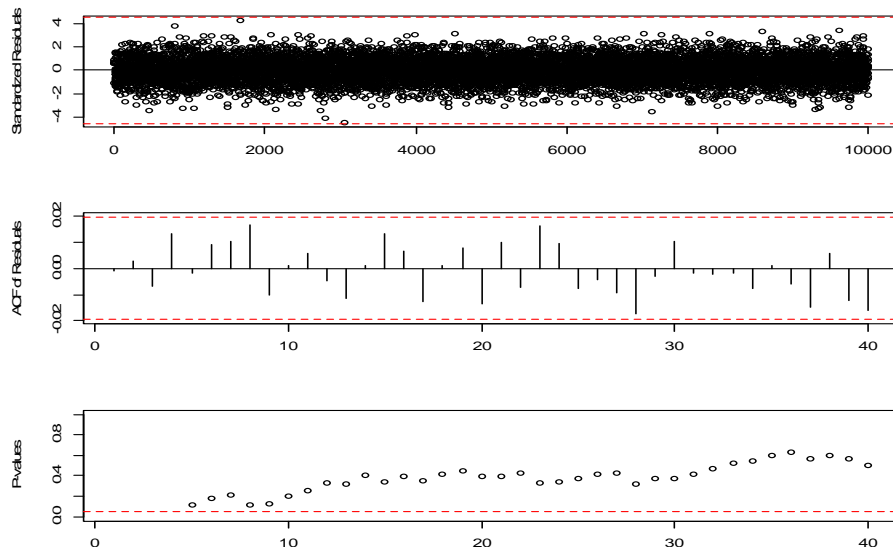


Fig. 7.8

Se vede că *acf* este suficient de mică pentru putea fi considerată zero, iar testul Ljung–Box indică faptul că în general pentru aproape toți h p-valoarea testului este

suficient de mare pentru a admite că rezidualii sunt iid. Acum dacă cercetăm modelul determinat găsim:

```
> model1

Call:
arima(x = s1, order = c(2, 0, 2), method = "ML")

Coefficients:
      ar1      ar2      ma1      ma2  intercept
    0.3831 -0.1976 -0.2773  0.5047     0.0016
s.e.  0.0289   0.0288   0.0258  0.0239     0.0152

sigma^2 estimated as 1.018:    log likelihood = -14276.47,    aic =
28562.93
```

Vedem că au fost estimați destul de bine coeficienții cu care s-a simulat seria ARMA.

d. Predicții

După validarea modelului de serie de timp acesta poate fi utilizat pentru predicția valorilor viitoare. În general pentru prezicerea valorilor unei serii temporale unde cunoaștem $(X_t)_{t=0,1,\dots,T}$ pentru valori $t > T$ definim predicția \hat{X}_{T+h} ca proiecția lui X_{T+h} pe spațiul generat de $(X_t)_{t=0,1,\dots,T}$. Această proiecție poate fi calculată ținând seama de produsele scalare conținute în funcția de autocorelație. Pentru seriile staționare calculul conduce la rezolvarea unui sistem asemănător cu (7.2.15) unde ϕ sunt necunoscute (sistemul Yule-Walker). Acest sistem poate fi rezolvat recursiv prin algoritmul Durbin-Levinson. O altă metodă de determinare a proiecției valabilă și pentru serii nestacionare este algoritmul inovațiilor. Detalii se pot găsi în cărțile citate ale lui Brockwell și Davis. În sistemul R există funcția *predict* care pe baza modelului (în general ARIMA) face predicția valorilor pentru un număr de pași temporali dați. Forma funcției este:

```
predict(object, n.ahead = 1, newxreg = NULL,
        se.fit = TRUE, ...)
```

În funcție de tipul de obiect R (*lm*, *ts*, *glm*, *loess*, *nls*, *poly*, *princomp*, *spline*) căruia i se aplică funcția *predict* poate avea argumente diferite în partea Pentru predicția seriilor ARIMA *object* din funcția *predict* trebuie să fie un obiect rezultat al funcției *arima*. Următorul script produce o predicție pentru 20 pași în avans pentru o serie ARMA(2,2) simulată ca în exemplul de precedent.

```
s1<-arima.sim(model=list(ma=c(-.3,.5), ar=c(0.4,-0.2)), n=1000)
model1<-arima(s1,order=c(2,0,2),method="ML")
#tsdiag(model1)
p1<-predict(model1, n.ahead=50)
p1
plot(s1, main="Seria s1 de 1000 masuratori si predictia pentru 50
pasi", xlab="Timpul")
lines(p1$pred, col=5,lwd=3)
```

Obținem următorul rezultat la consolă

```

$pred
Time Series:
Start = 1001
End = 1050
Frequency = 1
[1] 0.30741965 -0.01205071 -0.15271697 -0.10329277 -0.04444154 -0.03973509
[7] -0.05590263 -0.06270560 -0.06009920 -0.05718103 -0.05699457 -0.05781187
[13] -0.05814027 -0.05800335 -0.05785878 -0.05785189 -0.05789317 -0.05790899
[19] -0.05790182 -0.05789466 -0.05789444 -0.05789652 -0.05789728 -0.05789691
[25] -0.05789656 -0.05789655 -0.05789666 -0.05789669 -0.05789667 -0.05789666
[31] -0.05789666 -0.05789666 -0.05789666 -0.05789666 -0.05789666 -0.05789666
[37] -0.05789666 -0.05789666 -0.05789666 -0.05789666 -0.05789666 -0.05789666
[43] -0.05789666 -0.05789666 -0.05789666 -0.05789666 -0.05789666 -0.05789666
[49] -0.05789666 -0.05789666

$se
Time Series:
Start = 1001
End = 1050
Frequency = 1
[1] 0.9632667 0.9663791 1.0064598 1.0089120 1.0107742 1.0116303 1.0116409
[8] 1.0117388 1.0117442 1.0117491 1.0117511 1.0117512 1.0117514 1.0117514
[15] 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514
[22] 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514
[29] 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514
[36] 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514
[43] 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514 1.0117514
[50] 1.0117514

```

Obținem următorul grafic cu valorile seriei și valorile prezise:

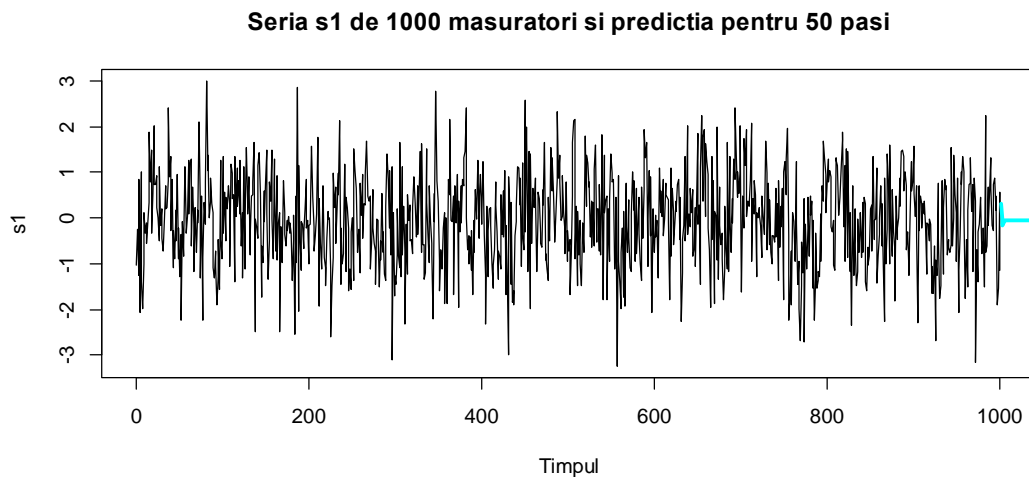


Fig. 7.9

Constatăm că predicțiile produse de model nu urmează variațiile seriei.

Aplicăm aceeași tehnică de tratare pentru seria reală de CO menționată la începutul secțiunii de serii de timp obținem:

```

>co=read.table(file.choose())
>tsco<-ts(co, frequency=168)
> plot(tsco, main="CO intre 1 ian.2008 si 30. iun. 2008",
+xlax="Timpul in saptamani")

```

Graficul seriei este în figura următoare:

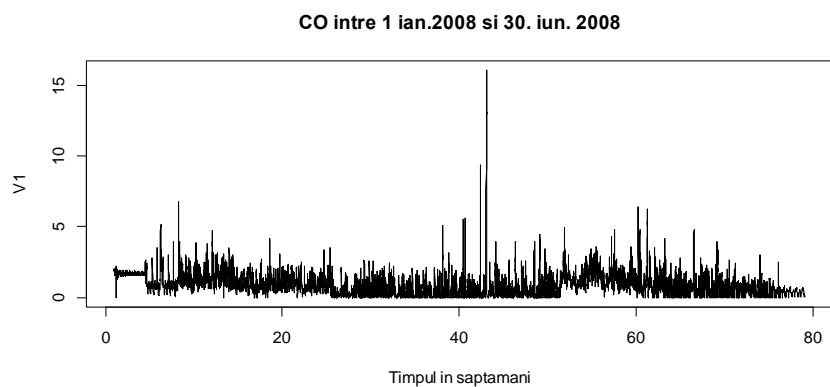


Fig. 7.10

Facem graficul pentru ACF și PACF pentru a vedea ce model ar fi potrivit.

```
> acf(tsco)
```

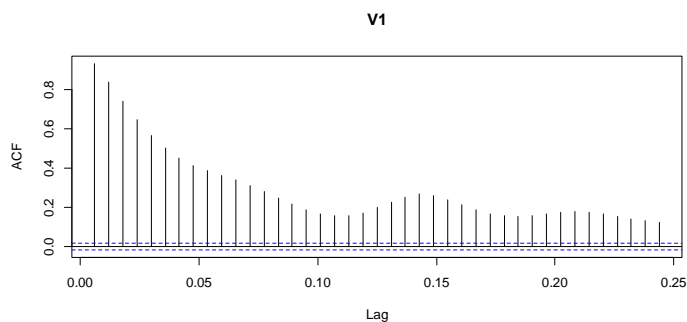


Fig. 7.11

Constatăm că ACF nu se anulează deci nu poate fi un model MA

Calculăm PACF:

```
> acf(tsco, type="partial")
```

Obținem:

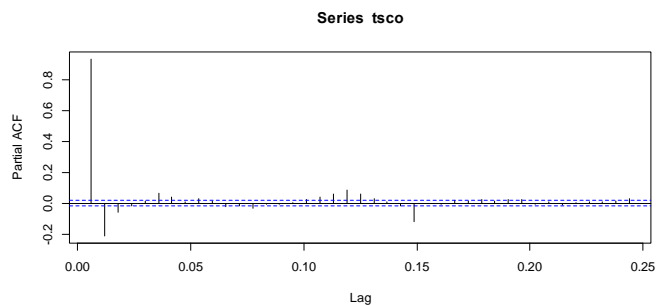


Fig. 7.12

Figura sugerează un model AR

Încercăm și cu EACF:

```
> eacf(tsko, ar.max = 20, ma.max = 20)
```

Obținem:

```
AR/MA
  0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
0 x x x x x x x x x x x x x x x x x x x x x
1 x x x o x x x x x x x x x x o x x x x x o
2 x o x o x x x x x x o x o x o o o o x x o
3 x o o x x x o o x x o x o o o o o o o o o
4 x x x x x x o x x o o o o o o o o o o o
5 x x x x o o o o x x o o x o o o o o o o
6 x x x x x o o o x x o o o o o o o o o o
7 x x x x x o o x x o x o x o o o o o o o
8 x x x x x o o x x o o o o o o o o o o o
9 x x x x x x x x x o o o x o o o o o o o
10 x x x x x x x x x o o o x o o o o o o o
11 x x x x x x x x x o o o x x o o o o o o
12 x x x x o x x x x x o o o o o o o o o o
13 x x x x o x x x x x o x o o o o o o o o
14 x x x x x x x x x x x x o o o o o o o
15 x x o x x o x x x x x x x x o o o o o
16 x x x x x x x x x x o o x x x o o o o o
17 x x o x x o x x o x x o o x x x o o o o
18 x x o x o o x x x x o o o x x o o o o x
19 x o x x o x o o o x x o o x x x x x o o
20 x x x x x x o o o o x x o x x x x x x o o
>
```

Tabloul obținut ar sugera un model ARMA(p,q) cu p în jur de 7-11 și q 9 sau 10. Dacă nu suntem prea pretențioși ar merge parcă p și q în jur de 5. Încercând am obținut mesajul că sunt probleme de convergență. Încercăm cu valori mai mici.

```
> m1<-arima(tsko,order=c(3,0,2),method="ML")
> tsdiag(m1)
```

Obținem:

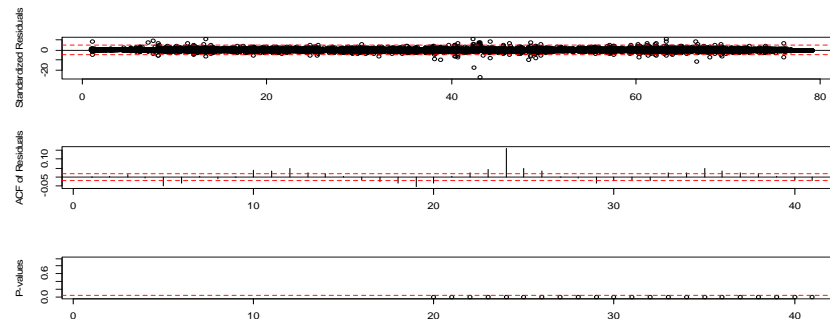


Fig. 7.13

ACF arată că există corelații semnificative între reziduali pentru diferențe mari de timp. Mergem totuși mai departe. Predicția dată de model arată astfel:

```
p1=predict(m1,n.ahead=168)
plot(p1$pred, col=5,lwd=3, main="Predictii CO")
```

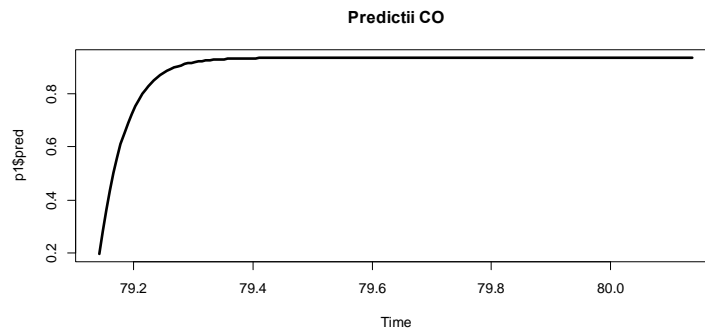


Fig. 7.14

Predicțiile pentru CO nu mai urmăresc variațiile constatate în modelul real. Se pare că nu am găsit modelul potrivit de serie temporală.

În general aplicate la modele reale teoria nu se potrivește întotdeauna cu realitatea. În secțiunea următoare prezentăm o altă tehnică de analiză a seriilor de timp care s-a dovedit mai bună pentru aceste date (CO).

7.2.2 Serii temporale deterministe și haos

Metoda sistemelor haotice, așa cum este descrisă în lucrarea,²⁴ permite într-o anumită măsură descrierea și predicția evoluției unor serii de timp. O implementare practică a acestor metode se găsește în pachetul gratuit de programe TISEAN disponibil publicului pe <http://www.mpipks-dresden.mpg.de/~tisean> (Institutul Max Planck pentru Fizică, Dresda). Autorii pachetului îl descriu într-o lucrare²⁵ disponibilă liber pe internet.

Ideea de bază a metodei este că datele urmează o evoluție deterministă

$$\bar{x}' = \bar{X}(\bar{x}) \quad (7.2.21)$$

cu \bar{x} într-un anumit spațiu normat finit sau infinit dimensional E , după o anumită perioadă de timp traiectoria sistemului dinamic (7.2.21) evoluează în jurul unui atractor A . Dinamica în A poate fi surprinsă efectiv de aproape orice funcție $s: E \rightarrow R$.

Fie $\Phi(t, x)$ fluxul lui (7.2.21) pentru o valoare fixată a lui $t = \tau$. Fie $F(\bar{x}) = \Phi(\tau, \bar{x})$, $F^{-1}(\bar{x}) = \Phi(-\tau, \bar{x})$, și fie $s: U \rightarrow R$ definită într-o vecinătate a lui A . Atunci, conform unei teoreme a lui F. Takens²⁶ în condiții generale pentru s și A și pentru un $m \in N$ suficient de mare, aplicația

$$S: A \rightarrow R^m$$

$$S(\bar{x}) = (s(F^{-(m-1)}(\bar{x})), s(F^{-(m-2)}(\bar{x})), \dots, s(F^{-1}(\bar{x})), \bar{x}) \quad (7.2.22)$$

este o scufundare a lui A în R^m .

Pe imaginea atractorului $S(A) \subset R^m$ avem

$$S(\Phi(n\tau, \bar{x}_0)) = \bar{s}_n = (s_{n-(m-1)}, s_{n-(m-2)}, \dots, s_{n-1}, s_n) \quad (7.2.23)$$

Dinamica lui F de pe atractorul A este exprimată pe imaginea sa $S(A)$ prin

$$\bar{s}_n \rightarrow G(\bar{s}_n) = \bar{s}_{n+1}$$

ca în figura următoare:

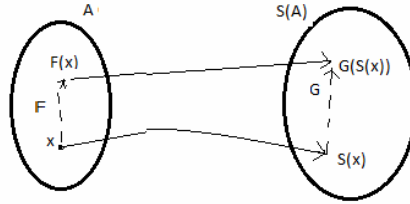


Fig. 7.15

O ușoară modificare a lui (7.2.233) introducând o anumită întârziere este $d \in N^*$ conduce la

$$\bar{s}_n = (s_{n-(m-1)d}, s_{n-(m-2)d}, \dots, s_{n-d}, s_n) \quad (7.2.24)$$

²⁴ Eckmann, J.P., D. Ruelle, 1985. *Ergodic theory of chaos and strange attractors*, Rev. Mod. Phys. 57, pp. 617-656

²⁵ Rainer Hegger, Holger Kantz, and Thomas Schreiber, 1999. *Practical implementation of nonlinear time series methods: The TISEAN package*, Chaos 9, pp. 413-435.

²⁶ H. D. I. Abarbanel, Reggie Brown, John J. Sidorowich, and Lev Sh. Tsimring, 1993. *The analysis of observed chaotic data in physical systems*, Rev M Phys, 65, pp. 1331-1392.

Intârzierea optima sugerată de Fraser și Swinney pentru (7.2.24) este recomandată ca fiind primul minim local al informației mutuale empirice:

$$H = \sum_{i,j} p_{i,j}(d) \ln \left(\frac{p_{i,j}(d)}{p_i p_j} \right) \quad (7.2.25)$$

Probabilitatile p_i și $p_{i,j}(d)$ în (7.2.25) sunt obținute prin partiționarea datelor $(s_n)_{1 \leq n \leq T}$ într-un număr de subintervale al valorilor datelor măsurate: p_i este probabilitatea ca funcția s să ia valori în subintervalul i iar $p_{i,j}(d)$ este probabilitatea de tranziție de la subintervalul i la subintervalul j după d unități.

Pentru a determina m , dimensiunea de scufundare, folosim metoda celui mai apropiat vecin fals (vezi referințele în subsolul paginii precedente). Ideea este că dacă \bar{s}_j este cel mai apropiat vecin al lui \bar{s}_i în R^m atunci $\bar{s}_{i+l} = G(\bar{s}_i)$ și $\bar{s}_{j+l} = G(\bar{s}_j)$ sunt de asemenea apropiați și raportul $R_i = \|\bar{s}_{i+l} - \bar{s}_{j+l}\| / \|\bar{s}_i - \bar{s}_j\|$ este mai mic decât un prag euristic R_t . În caz contrar punctul \bar{s}_j este marcat ca vecin fals. Dacă procentul de puncte cu un cel mai apropiat vecin fals este prea mare, considerăm ca m este prea mică.

Metode de predicție

Predicția de ordin zero este data de

$$s_{n+k} = \frac{1}{\text{card}(U_n)} \sum_{\bar{s}_j \in U_n} s_{j+k} \quad (7.2.26)$$

unde U_n este o mică vecinătate a lui \bar{s}_n în R^m . Predicția este făcută pentru $n+k > T$, ultima măsurătoare.

Predicția local liniară este dată de formula

$$s_{n+l} = \bar{a}_n \bar{s}_n + b_n \quad (7.2.27)$$

unde coeficienții $\bar{a}_n \in R^m$, $b_n \in R$ sunt determinați astfel încât

$$\sigma^2 = \sum_{\bar{s}_j \in U_n} (s_{j+l} - \bar{a}_n \bar{s}_j - b_n)^2 \quad (7.2.28)$$

să fie minimă. Suma este efectuată după toți \bar{s}_j dintr-o mică vecinătate U_n a lui \bar{s}_n dar destul de mare încât să conțină suficiente puncte pentru a asigura ca problema celor mai mici patrate (7.2.28) să fie nesingulară. În lucrările citate în pagina anterioară se găsesc și alte metode de predicție.

Așa cum am mai menționat aceste metode au fost implementate de Rainer Hegger, Holger Kantz și Thomas Schreiber în pachetul TISEAN care este disponibil și în R. Pentru a putea fi folosit în R trebuie ca pachetul TISEAN cu executabilele să fie dezarhivat într-un director. Pe urmă trebuie instalat în R pachetul RTisean care conține apeluri în stilul R pentru executabilele din TISEAN. La primul apel al unei funcții din RTisean se cere directorul cu executabilele TISEAN. Informații despre funcțiile din pachet pot fi obținute prin `?RTisean`. Acest pachet de analiză a seriilor temporale a fost utilizat pe date de poluare obținute în municipiul București, în

particular pentru concentrația de CO al cărei grafic este în Fig. 7.1. Reproducem²⁷ câteva rezultate obținute în predicția poluării utilizând tehnici din această secțiune. Calculele au fost făcute cu rutinele TISEAN apelate direct (nu din R, dar o examinare a fișierului de help din RTisean arată că apelul din R se face cu aceeași parametri ca și apelul direct din consola windows). Rezultatele au fost salvate în fișiere text de unde au fost trecute în Excel pentru examinare și pentru reprezentare grafică. Utilizând RTisean nu mai este nevoie de acest șir de operații care încetinește mult analiza de variante. Mai jos reproducem câteva grafice unde sunt reprezentate valorile măsurate și cele prezise pe baza măsurătorilor anterioare. Predicțiile sunt pentru o săptămână (168 ore).

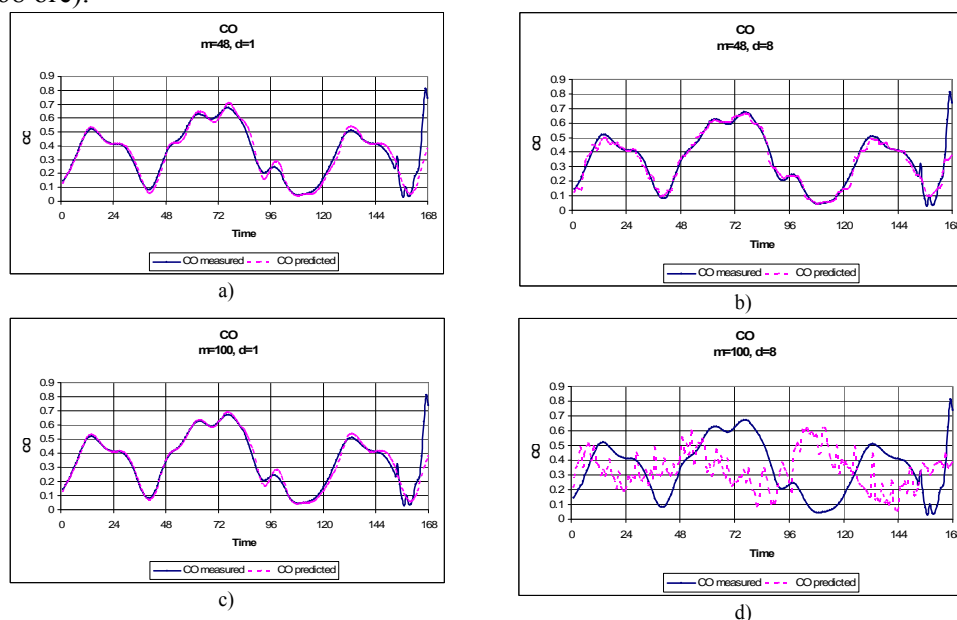


Fig. 7.16

Concentrația de CO este nefiresc de bine estimată prin predicția 7.2.26 pentru o dimensiune de scufundare $m \geq 48$. Contrar cu teoria, la o întârziere $d=1$ predicția e mai bună decât pentru $d=8$, unde este primul minim al informației mutuale empirice 7.2.25. Rutina TISEAN utilizată a fost lzo (în R se numește lzo.test).

Nu au fost la fel de bune predicțiile pentru alți poluanți. Pentru particulele PM 2.5 se vede din graficele de mai jos.

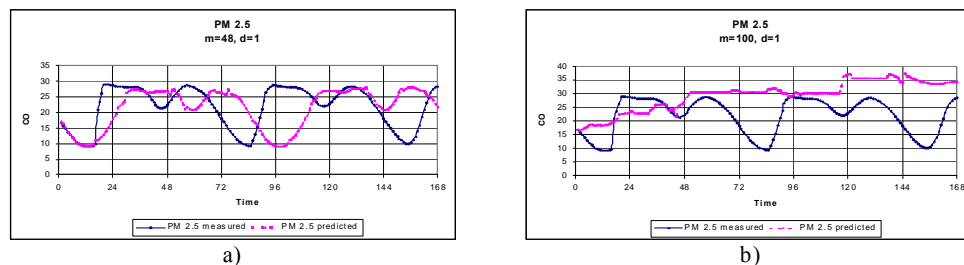


Fig. 7.17

Predicțiile pentru 168 de ore a poluării cu PM 2.5

²⁷Viorel Petrehus, Ileana Armeanu, Camelia Slave, *Analysis of the Urban Air Pollution in Bucharest*, Proceedings BALCOR 2011, vol 2. pp 72-78

În momentul când s-au făcut predicții simultane (funcția s din (7.2.22) are valori în R^p , deci S are valori în R^{mp} , pentru p poluanți considerați simultan) am constatat prin multe încercări că predicțiile nu devin mai bune. În figura următoare vedem ce se întâmplă dacă simultan luăm CO și PM 2.5 în considerare. Datele măsurate pentru un an și jumătate sunt utilizate pentru a prezice poluarea pe o săptămână.

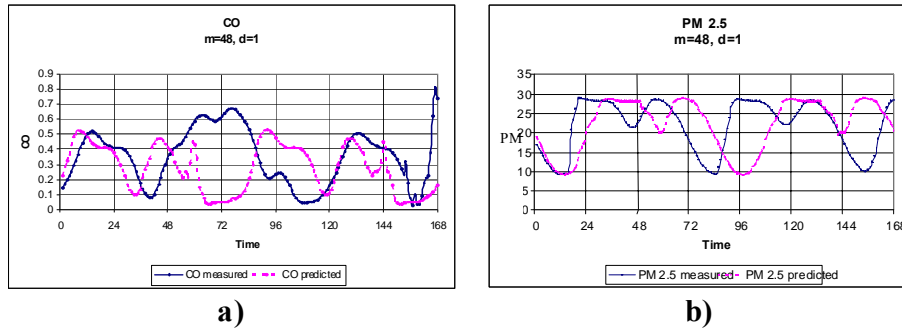


Fig. 7.18

Cele mai bune predicții simultane sunt obținute pentru $d = 1$ și sunt semnificative numai pentru un interval de timp mai mic de 20 de ore

Din experimentele numerice am ajuns la următoarele concluzii:

1. Întârzierea $d = 1$ este o alegere bună, uneori mai bună ca întârzierea la care informația mutuală empirică ajunge la un minim local.
2. O alegere bună pentru vecinătatea U_n , care determină câte puncte intră în calcul în formula (7.2.26) se află experimental. Pentru predicția local liniară o alegere bună a numărului de puncte în formula (7.2.28) este de cel puțin patru ori numărul de parametri de estimat.
3. Alegerea dimensiunii de scufundare astfel încât fracțiunea de puncte având cel mai apropiat vecin fals este mai mică de 5%, dă rezultate bune. De asemenea, este indicat să alegem m , astfel ca vectorul \bar{s}_n să cuprindă datele unei perioade în cazul în care se suspectează că datele au o anumită periodicitate.
4. Este posibil ca datele să nu respecte o lege deterministă. În acest caz, orice alegere a lui m și d nu este o alegere bună.
5. Previziunile pe termen scurt sunt în general mai bune ca cele pe termen lung.
6. Prin experimentare numerică ajungem la un moment dat să găsim o dimensiune m de scufundare și o întârziere d care utilizate pentru anumite date dintr-un anumit loc duc la previziuni acceptabile. Este de așteptat ca acestea să reflecte particularități ale sistemului dinamic (7.2.21) din acel loc și să fie valabile și pentru alt set de măsurători. Acest fenomen, determinarea parametrilor de scufundare m și d , ar fi analogul determinării parametrilor într-un model statistic (de exemplu ARMA(p, q)).

Bibliografie

** The R Development Core Team, *R: A Language and Environment for Statistical Computing*, 2011-12-13.

1. M. Allehard. *A Tyny Handbook of R*, Springer, 2011-12-13.
2. Y. Aragon. *Séries Temporelles avec R. Méthodes et cas*. Springer, 2011.
3. M. Dumitrescu, A. Bătătorescu. *Applied Statistic using the R System*, Ed. Universității București, 2006.
4. J. D. Cryer, Kung-Sik Chan. *Time Series Analysis, with Applications in R*, Springer, 2008.
5. J. Maindonald, W. J. Braun. *Data Analysis and Graphics using R – an Example-Based Approach*, Cambridge University Press, 2003
6. Murrell, P. (2005) *R Graphics*. Chapman & Hall/CRC Press.
7. A. Păun, M. Păun. *Analiză Statistică folosind limbajul R*, Editura Matrix Rom, București, 2009.
8. V. Petrehuș, S.A. Popescu, „Probabilități și Statistică”, Universitatea Tehnică de Construcții, București, 1997, civile.utcb.ro/cmat/cursrt/psvp.pdf.
9. D. Ruppert. *Statistics and Data Analysis for Financial Engineering. Use R!* Springer, 2010.
10. G. Zbăganu. *Metode Matematice în Teoria Riscului și Actuariat*, Editura Universității București, 2004.