

Analiza datelor în R

Curs 2

Pachete

- ▶ Pachete: colecții de funcții, documentație aferentă și seturi de date
- ▶ În prezent sunt disponibile >13000 pachete R
(https://cran.r-project.org/web/packages/available_packages_by_name.html).
- ▶ Instalare:
`install.packages("nume_pachet")`
- ▶ Încărcare în sesiunea de lucru curentă:
`library(nume_pachet)`
- ▶ Identificarea pachetelor încărcate pentru sesiunea de lucru curentă:
`sessionInfo()`
- ▶ Un set de date dintr-un pachet deja încărcat se aduce în workspace cu
`data(nume_set_date)`

Funcții pentru vectori

1. Funcții specifice pentru vectori logici

- ▶ `all(x)` → T dacă toate componentele lui x sunt T; F în rest
- ▶ `any(x)` → T dacă măcar o componentă a lui x e T; F în rest
- ▶ `which(x)` → pozițiile componentelor T din x

2. Funcții pentru vectori numerici

- ▶ `min(x)`, `max(x)`
- ▶ `range(x)` → cea mai mică și cea mai mare valoare din x
- ▶ `which.min(x)`, `which.max(x)` → poziția primei apariții în x pentru cea mai mică, respectiv cea mai mare valoare

- ▶ `sum(x)` → suma elementelor lui `x`
- ▶ `prod(x)` → produsul elementelor lui `x`
- ▶ `cumsum(x)` → vectorul sumelor cumulate
- ▶ `cumprod(x)` → vectorul produselor cumulate
- ▶ `diff(x)` → vectorul diferențelor între elemente succesive

3. Funcții generale

- ▶ `length(x)` → numărul de elemente din `x`
- ▶ `rev(x)` → vectorul elementelor în ordine inversă
- ▶ `sort(x)`, `sort(x, decreasing=T)` - funcții de ordonare
- ▶ `order(x)` → pozițiile elementelor lui `x` ordonate crescător după valoare
- ▶ `unique(x)` → elimină valorile care se repetă din `x`

Simbolul NA

NA (not available) - codificare pentru date lipsă

Orice operație care implică NA are ca rezultat NA.

- ▶ Unele funcții au opțiunea `na.rm=TRUE`, care elimină valorile lipsă înainte de efectuarea calculelor.
- ▶ La sortare crescătoare, valorile NA sunt puse pe ultimele poziții.
- ▶ Funcția `is.na(x)` returnează un vector de valori logice, respectiv T acolo unde componenta corespunzătoare a lui x este NA și F în rest.

Factori

Tipuri de date statistice:

- ▶ cantitative
 - ▶ discrete
 - ▶ continue

- ▶ calitative (catoriale)
 - ▶ binare
 - ▶ nominale
 - ▶ ordinale

Reprezentare în R:

numeric

factor

O structură factor reprezintă valorile (categoriile) variabilei calitative printr-un vector de numere întregi $\in \{1, 2, \dots, k\}$ (k = numărul de categorii) și un vector intern de șiruri de caractere corespunzătoare valorilor întregi.

Factori

Exemplu:

```
optiuniVot=c("C","A","B","C","C","B","A")  
class(optiuniVot)
```

→ *character*

```
optiuniVot=factor(optiuniVot)  
class(optiuniVot)
```

→ *factor*

```
levels(optiuniVot)
```

→ *A, B, C*

Variabila `optiuniVot` este reprezentată intern prin vectorul 3, 1, 2, 3, 3, 2, 1, cu convenția 1=A, 2=B, 3=C. Implicit, nivelele sunt codificate în ordinea crescătoare a tipului de date din care factorul a fost construit.

Se poate opta pentru o altă codificare a nivelelor decât cea implicită:

```
optiuniVot=factor(optiuniVot, levels=c("B",  
                                         "A", "C"))  
levels(optiuniVot)
```

→ *B, A, C*

```
as.numeric(optiuniVot)
```

→ 3, 2, 1, 3, 3, 1, 2

Factori ordonați

Pentru reprezentarea variabilelor calitative ordinale, acolo unde relația de ordine între categorii este importantă, se utilizează la crearea variabilei factor opțiunea `ordered=T` și se specifică nivelele în ordinea crescătoare dorită.

Exemplu:

```
x=c("mic", "mediu", "mic", "mare", "mediu", "mic")  
x=factor(x, ordered=T, levels=c("mic", "mediu", "mare"))
```

Două variabile de tip factor ordonat care au aceleași nivele și aceeași structură de ordine pot fi comparate (element cu element).

Dataframe-uri

Cea mai uzuală modalitate de reprezentare a seturilor de date statistice este cea de tabel, unde

- ▶ liniile = obiecte (indivizi)
- ▶ coloanele = attribute (variabile)

→ În R: **dataframe**

Coloanele=vectori de aceeași lungime. Ele pot avea tipuri diferite.

Un dataframe poate fi creat:

- ▶ dintr-un set de vectori preexistenți

```
nume=c("Popescu", "Ionescu", "Anton")
prenume=c("Mircea", "Alina", "Andrei")
varsta=c(24, 31, 45)
bd=data.frame(nume, prenume, varsta)
```

- ▶ dintr-un set de date extern

Importarea / exportarea seturilor de date

Import

- ▶ date în format .csv:

```
numeDF=read.csv("numefisier.csv") (! opțiuni)
```

- ▶ date în format .txt:

```
numeDF=read.table("numefisier.txt")
```

- ▶ alte formate proprietare → pachetul *foreign*

Export

```
write.csv(numeDF, "numefisier.csv")
```

Pentru detalii, a se vedea R Data Import/Export manual

<https://cran.r-project.org/doc/manuals/R-data.pdf>

Funcții pentru dataframe-uri

- ▶ `attach(umeDF)` → atașarea unui dataframe la sesiunea de lucru curentă
- ▶ `detach(umeDF)` → eliminarea dataframe-ului din sesiunea curentă
- ▶ `str(umeDF)` → informații despre conținut
- ▶ `names(umeDF)` → afișarea numelor coloanelor (variabilelor)
- ▶ `dim(umeDF)` → numărul de linii și de coloane
- ▶ `edit(umeDF)` → deschidere în format spreadsheet editabil pentru modificarea conținutului

Exemplu: Setul de date *Cars93* din pachetul *MASS*

Selectarea datelor dintr-un dataframe

- ▶ O variabilă `v` dintr-un dataframe `d` se accesează cu `d$v`.
Dacă dataframe-ul a fost deja atașat la sesiunea de lucru curentă, variabilele sale pot fi accesate direct utilizând numele lor.
- ▶ selectarea liniei `i`: `d[i,]`
- ▶ selectarea coloanei `j`: `d[, j]`
- ▶ selectarea elementului de pe poziția `(i,j)`: `d[i, j]`
- ▶ selectarea mai multor linii (coloane) specificate:
`d[istart:ifinal,]`
`d[c(i1,i2,...,ik),]`
`d[-c(i1,i2,...,ik),]`
`d[, jstart:jfinal]`
`d[, c(j1,j2,...,jk)]`
`d[, -c(j1,j2,...,jk)]`
- ▶ selectare condiționată a obiectelor:
`d1=subset(d, conditie)`

Exerciții

Se consideră setul de date *Cars93* din pachetul MASS.

- i) Să se construiască un dataframe care conține doar variabilele Manufacturer, Make, Price, Passengers și Origin.
- ii) Să se selecteze și afișeze doar datele pentru mașinile de proveniență americană. Să se determine câte astfel de mașini sunt în baza de date.
- iii) Să se afișeze modelele de mașini al căror producător este Ford.
- iv) Să se selecteze și afișeze datele pentru mașinile care pot transporta cel puțin 5 pasageri, sortate crescător după preț.