

Analiza datelor în R

Curs 6

Analiză exploratorie multivariată

Scop: descrierea relațiilor dintre variabile

- ▶ cantitative vs. calitative
- ▶ calitative vs. calitative
- ▶ cantitative vs. cantitative

Relația dintre o variabilă cantitativă și una calitativă

1. Vom considera setul de date *iris* din *datasets*.

Vom analiza lungimea petalelor, lățimea petalelor, lungimea se-palelor și lățimea se-palelor, comparativ pe specii:

- ▶ calcul de indicatori numerici
- ▶ reprezentare grafică (boxplot grupat pe categorii):
`boxplot (varCantitativa~varCalitativa)`

Relația dintre o variabilă cantitativă și una calitativă

2. Vom utiliza setul de date *energy* din pachetul *ISwR* pentru a compara consumul de energie (*expend*) pentru 22 femei, în funcție de clasa de greutate corporală (*stature* = *lean* sau *obese*).

- ▶ boxplot grupat
- ▶ Atunci când grupurile analizate sunt mici, este preferabil să se reprezinte grafic datele sub formă de puncte (stripchart):

```
stripchart(expend~stature,method="jitter")
```

Relația dintre două sau mai multe variabile calitative

Se utilizează pentru reprezentare tabele de contingență.

Introducerea unui tabel preexistent

Se analizează preferința pentru un anumit tip de băutură (cola, cafea sau ceai) pe un eșantion de 105 persoane și se obțin răspunsurile din tabelul de mai jos:

		băutura		
		cola	cafea	ceai
F/M	F	10	12	28
	M	20	31	4

Datele se preiau într-o matrice, care se convertește în tabel cu funcția `as.table`.

Relația dintre două sau mai multe variabile calitative

Cod R:

```
bauturi=matrix(c(10,12,28,20,31,4),byrow=T,nrow=2)
```

```
# denumirea liniilor si coloanelor  
colnames(bauturi)=c("cola","cafea","ceai")  
rownames(bauturi)=c("F","M")  
bauturi
```

```
# denumirea variabilelor reprezentate pe linii  
# si coloane  
names(dimnames(bauturi))=c("F/M","bautura")  
bauturi
```

```
# conversie in tabel  
bauturi=as.table(bauturi)  
bauturi
```

Relația dintre două sau mai multe variabile calitative

Frecvențe marginale=frecvențe pentru fiecare variabilă calitativă

`margin.table(bauturi,1)` → frecvențe marginale pentru atributul de pe linii

`margin.table(bauturi,2)` → frecvențe marginale pentru atributul de pe coloane

Tabele de frecvențe relative

`prop.table(bauturi,1)` → proporții pentru fiecare celulă din totalul de pe linie (frecvențe relative pentru fiecare nivel al atributului de pe linii)

`prop.table(bauturi,2)` → proporții pentru fiecare celulă din totalul de pe coloană

`prop.table(bauturi)` → proporții pentru fiecare celulă din totalul pe tabel

Relația dintre două sau mai multe variabile calitative

Tabularea datelor dintr-un dataframe: funcțiile `table` și `xtabs`

Setul de date *UCBAdmissions* conține informații despre situația admiterilor la Universitatea Berkeley în 1973.

```
xtabs(Freq~Gender+Admit)
```

```
prop.table(xtabs(Freq~Gender+Admit), 1) → proporția  
de admiși pe sexe
```

Se observă că, dintre aplicanții bărbați, 44.5% au fost admiși, în timp ce dintre femei, doar 30.4%. Analiza situației pe fiecare departament arată însă că procentajele de femei și bărbați admiși sunt comparabile.

→ **paradoxul lui Simpson:** un trend care apare într-o combinație de grupuri dispare sau se inversează în interiorul fiecărui grup.

Relația dintre două variabile cantitative

Se analizează pe un eșantion de volum n două caracteristici, X și Y , și se obțin perechile de valori $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.

- ▶ grafic de împrăștiere: `plot(x, y)`
- ▶ coeficientul de corelație liniară (Pearson)

$$r_{XY} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- ▶ $r \in [-1, 1]$
- ▶ $r < 0$ - corelație liniară negativă
- ▶ $r \approx 0$ - nu există asociere liniară între variabile
- ▶ $r > 0$ - corelație liniară pozitivă
- ▶ r este invariant la translație sau scalare, deci datele pot fi standardizate.

În R: `cor(x, y)`

Relația dintre două variabile cantitative

1. Analizăm relația dintre înălțime și greutate pentru 15 femei cu vârstele între 30 și 39 de ani (setul de date *women* din *datasets*).
2. Setul de date *mtcars* din *datasets* conține caracteristicile tehnice a 32 modele de mașini din perioada 1973-1974. Analizăm dacă există o relație între greutatea unei mașini (*wt*) și consumul de combustibil (*mpg*).

Dacă relația dintre variabile nu este liniară dar este monotonă, se poate calcula coeficientul de corelație Spearman ρ = coeficientul Pearson calculat pentru rangurile observațiilor. Coeficientul Spearman poate fi calculat și pentru date ordinale.

În R: `cor(x, y, method="spearman")`

Relația dintre două variabile cantitative

Cantitățile care se modifică multiplicativ (i.e., rata de schimbare este proporțională cu cantitatea) pot fi logaritmate înainte de a fi analizate.

Exemplu:

Setul de date *Animals* din *datasets* conține informații despre greutatea corporală (kg) și greutatea creierului (g) la mai multe specii de animale. Vom analiza relația dintre cele două cantități în forma dată, apoi logaritmuate.

Grafice quantilă - quantilă

În cazul a două serii de date (i.e. măsurători ale aceleiași caracteristici în două grupuri distincte), se utilizează pentru a determina dacă acestea urmează aceeași distribuție.

- ▶ grupuri de aceeași mărime: se ordonează crescător valorile în cele două grupuri, i.e. $x[1] \leq x[2] \leq \dots x[n]$, $y[1] \leq y[2] \leq \dots y[n]$, apoi se reprezintă grafic punctele $(x[i], y[i])$
- ▶ grupuri de mărimi diferite: se reduce grupul mai mare la dimensiunea celui mai mic păstrând *Min*, *Max* și alegând quantile echidistante între acestea.

Dacă distribuțiile sunt similare, punctele se vor găsi aproximativ pe prima bisectoare.

În R: `qqplot(x, y)`

Exerciții

1. Pentru setul de date *airquality* din *datasets*, să se reprezinte temperaturile înregistrate sub formă de boxplot-uri grupate pe luni.
2. Setul de date *reaction.time* din *UsingR* conține timpii de reacție la un stimul extern pentru participanții la un studiu, precum și informații despre vârsta și sexul acestora și dacă foloseau sau nu telefonul mobil la momentul respectiv. Să se reprezinte timpul de reacție în funcție de:
 - a) categoria de vârstă
 - b) sex
 - c) utilizarea sau nu a telefonului mobil.

Exerciții

3. În setul de date *twins* din *UsingR* se găsesc IQ-urile a 27 perechi de gemeni separați la naștere, dintre care unul a fost crescut de familia biologică (*Biological*), iar celălalt de o familie adoptivă (*Foster*). Să se determine dacă există o asociere între IQ-urile celor două categorii de gemeni.
4. Setul de date *kid.weights* din *UsingR* conține înălțimile și greutatea a 250 copii cu vârstele cuprinse între 0 și 12 ani. Să se reprezinte greutatea în funcție de înălțime și să se determine tipul relației între aceste două mărimi și coeficientul de corelație adecvat.

Exerciții

5. Setul de date *UNLifeExpectancy*, disponibil la adresa

<http://instruction.bus.wisc.edu/jfreese/jfreesebooks/RegressionModeling/BookWebDec2010/CSVData/UNLifeExp.csv>, conține informații cu privire la speranța de viață și alți parametri socio-economici ai țărilor lumii.

- a) Să se reprezinte histograma și curba de densitate pentru speranța de viață (*LIFEEXP*).
- b) Să se reprezinte histograma pentru cheltuielile de sănătate (*HEALTHEXPEND*).
- c) Să se determine dacă există o asocierie între speranța de viață într-o țară și cheltuielile cu sănătatea, respectiv speranța de viață și fertilitatea (*FERTILITY*) în țara respectivă.