

Analiza datelor în R

Curs 5

Indicatori statistici numerici

Se analizează o caracteristică de tip cantitativ pe un eșantion de volum n dintr-o populație și se obțin valorile x_1, x_2, \dots, x_n .

Indicatori de poziție

- Media de selecție

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

În R: `mean(x)`

- Mediana

Dacă presupunem seria de date ordonată crescător

$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$, atunci

$$Me = \begin{cases} x_{\frac{n+1}{2}}, & \text{dacă } n = 2k + 1 \\ \frac{x_{\frac{n}{2}} + x_{\frac{n}{2}+1}}{2}, & \text{dacă } n = 2k \end{cases}$$

În R: `median(x)`

Indicatori statistici numerici

- ▶ Valoarea modală = valoarea (valorile) cu cea mai mare frecvență în setul de date.

În R: `pachetul modeest`

- ▶ p -quantilele ($p \in (0, 1)$)
 $y = Q(p) \iff$ o proporție p dintre valorile x_1, \dots, x_n sunt mai mici decât y , iar restul sunt mai mari decât y

În R: `quantile(x, p)`

- ▶ quartile: $Q(0.25)$, $Q(0.5)$, $Q(0.75)$

În R: `summary(x)`

Indicatori statistici numerici

Indicatori ai împrăştierii

- ▶ Amplitudinea = $Max - Min$
- ▶ Lungimea intervalului interquartilic
 $IQR = Q(0.75) - Q(0.25)$
- ▶ Dispersia de selecție

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

În R: `var(x)`

- ▶ deviația (abaterea) standard de selecție $s = \sqrt{s^2}$

În R: `sd(x)`

Indicatori statistici numerici

Indicatori de formă a distribuției

- Asimetria ("skewness") =
$$\frac{\sum_{i=1}^n (x_i - \bar{x})^3}{ns^3}$$
- $< 0 \rightarrow$ asimetrie stânga
 - $\approx 0 \rightarrow$ distribuție simetrică
 - $> 0 \rightarrow$ asimetrie dreapta

În R: `skewness` din pachetul `moments`

- Aplatizarea ("kurtosis") =
$$\frac{\sum_{i=1}^n (x_i - \bar{x})^4}{ns^4}$$
- $< 3 \rightarrow$ distribuție platikurtică
 - $\approx 3 \rightarrow$ distribuție normokurtică
 - $> 3 \rightarrow$ distribuție leptokurtică

În R: `kurtosis` din pachetul `moments`

Exemplu

Setul de date *cfb* din pachetul *UsingR* conține informații dintr-un studiu asupra finanțelor consumatorilor efectuat de USA Federal Reserve. Vom ilustra conceptele de mai sus pentru analiza variabilelor *AGE* (vârsta participanților la studiu) și *INCOME* (veniturile acestora în 2001).

Tabele și grafice de frecvențe

Se folosesc pentru analizarea datelor categoriale sau a celor cantitative grupate pe categorii.

Tabele de frecvențe

- ▶ `table(x)` → valorile distincte din x și frecvențele lor absolute
- ▶ `table(cut(x, k))` → intervalul de valori ale lui x se împarte în k subintervale de lungime egală și se determină frecvența absolută pentru fiecare subinterval

Grafice de frecvențe

- ▶ `barplot(vectFrecv, names.arg=vectCateg)` sau `barplot(table(x))`
- ▶ `pie(vectFrecv, vectCateg)` sau `pie(table(x))`

Exemple

1. Setul de date *central.park.cloud* din pachetul *UsingR* conține informații cu privire la vremea din Central Park în luna mai 2003 (cer senin / parțial noros / noros). Să se construiască tabelul de frecvențe și graficele de frecvențe.
2. Setul de date *airquality* din *datasets* conține măsurători zilnice asupra calității aerului în New York în perioada mai - septembrie 1973. Să se construiască un tabel de frecvențe pentru variabila *Temp* și graficele de frecvențe corespunzătoare.

Histograme

Criterii de alegere a numărului de clase de grupare:

- ▶ la alegere
- ▶ $k \approx \sqrt{n}$ (Excel)
- ▶ $k = \lceil \log_2 n + 1 \rceil$ (Sturges)
- ▶ $k \approx 3.49 \cdot s \cdot n^{-1/3}$ (Scott)

În R: `hist(x)`

Argumente suplimentare:

- ▶ `breaks=m` $\rightarrow \approx m$ clase de grupare egale
- ▶ `breaks="Sturges"` (implicit) sau `"Scott"`
- ▶ `breaks=vect` unde *vect* este vectorul extremităților de intervale de grupare
- ▶ `prob=T` \rightarrow aria histogramei = 1

Curbe de densitate

Poligonul de frecvențe

- ▶ linie poligonală care unește mijloacele segmentelor orizontale superioare ale dreptunghiurilor histogramei

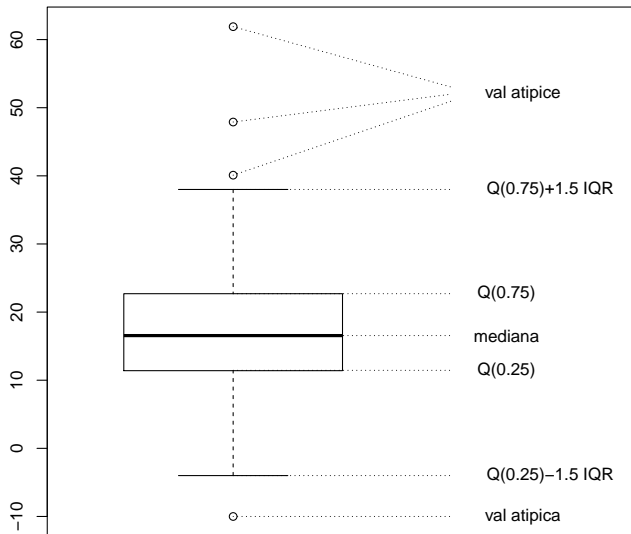
Curba de densitate

- ▶ versiune "continuă" a histogramei
- ▶ înălțimea curbei de densitate într-un punct estimează proporția de valori din eșantion care se găsesc într-un interval de lățime specificată ("bandwidth"), centrat în acel punct

În R: `plot(density(x))`

Exemplu: Histograma și curba de densitate pentru variabila *Temp* din *airquality*

Grafice de tip boxplot



În R: `boxplot(x)`

Grafice de tip quantilă - quantilă

- ▶ Se utilizează în cazul unui singur eșantion pentru a compara distribuția datelor cu o distribuție de referință (e.g. normală).
- ▶ se reprezintă grafic puncte având ca și coordonate:
 - ▶ pe Oy, valorile din eșantion, sortate crescător, fiecare reprezentând o anumită p -quantilă a respectivului eșantion
 - ▶ pe Ox, p -quantilele corespunzătoare ale distribuției teoretice de referință
- ▶ Dacă există concordanță, punctele se vor afla aproximativ pe o dreaptă.

În R: `qqnorm` (comparare cu distribuția normală); `qqplot`

Grafice de tip quantilă - quantilă

Exemplu:

```
x=rnorm(1000); y=rchisq(1000,df=5)
```

```
qqnorm(x)
```

```
qqnorm(y)
```

