

Analiza datelor în R

Curs 7

Statistici de selecție

Fie X_1, X_2, \dots, X_n variabile aleatoare independente, identic distribuite (i.i.d.). Definim:

- ▶ media de selecție $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$
- ▶ dispersia de selecție $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

Principii generale

- ▶ Se formulează cu privire la parametrul sau parametrii de interes ai populației o ipoteză nulă H_0 și o ipoteză alternativă H_a ;
- ▶ Pe baza datelor din eșantion se poate lua decizia nerespingerii lui H_0 sau a respingerii lui H_0 , i.e., a acceptării lui H_a .
- ▶ $P(\text{resping } H_0 \mid H_0 \text{ e adevărată}) = \alpha = \text{nivelul de semnificație al testului}$
- ▶ $P(\text{nu resping } H_0 \mid H_0 \text{ e falsă}) = \beta$
- ▶ $1 - \beta = P(\text{accept } H_a \mid H_a \text{ e adevărată}) = \text{puterea testului}$
- ▶ R. Fisher - Lady tasting tea experiment

Principii generale

- ▶ Se alege o statistică test a cărei distribuție este cunoscută când H_0 este adevărată și se calculează valoarea ei pe datele din eșantion;
- ▶ Dacă probabilitatea de a obține valori cel puțin "la fel de extreme" sub ipoteza nulă (**p-valoarea testului**) este foarte mică (mai mică decât α), vom respinge H_0 în favoarea lui H_a . În caz contrar, nu respingem H_0 .

Testul t - un eșantion

Este un test pentru media unei populații cu volum de eșantion mare ($n > 30$) sau distribuție (aproximativ) normală și volum de eșantion mic, cu dispersie necunoscută.

- ▶ $H_0 : \mu = \mu_0$
- ▶ $H_a : \mu \neq \mu_0$ sau $H_a : \mu < \mu_0$ sau $H_a : \mu > \mu_0$

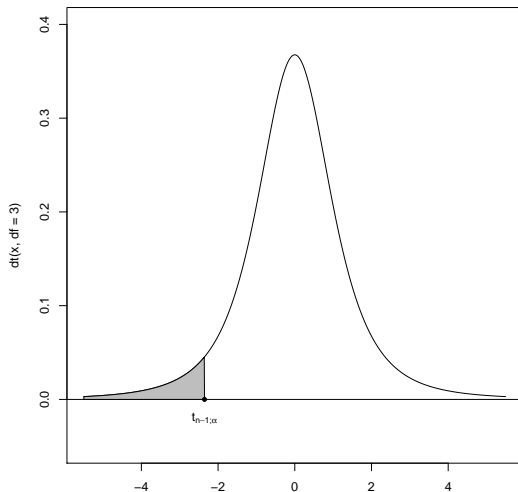
- ▶ Statistica test: $T = \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}}$

În ipoteza că H_0 e adevărată, $T \sim t(n - 1)$.

- ▶ $t_{obs} = \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}}$
- ▶ Dacă p-valoarea corespunzătoare lui t_{obs} este $< \alpha$, se respinge H_0 ; în caz contrar, nu se respinge H_0 la nivelul de semnificație ales.

Testul t - un eșantion

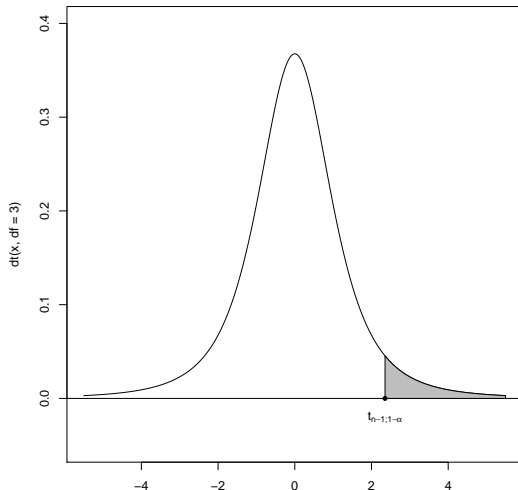
$$H_a : \mu < \mu_0$$



Regiunea critică: $(-\infty, t_{n-1; \alpha})$

Testul t - un eșantion

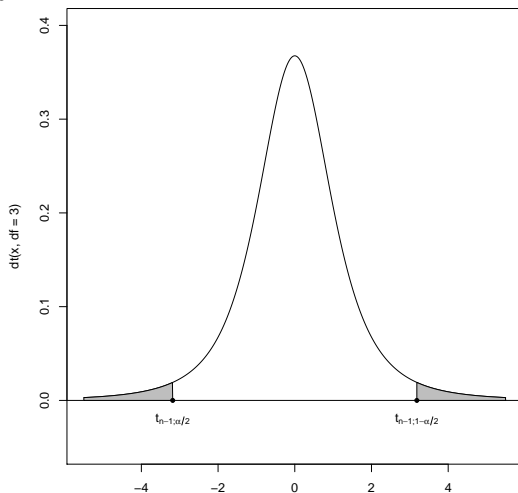
$$H_a : \mu > \mu_0$$



Regiunea critică: $(t_{n-1; 1-\alpha}, \infty)$

Testul t - un eșantion

$$H_a : \mu \neq \mu_0$$



Regiunea critică: $(-\infty, t_{n-1; \frac{\alpha}{2}}) \cup (t_{n-1; 1-\frac{\alpha}{2}}, \infty)$

Testul t - un eșantion

În R: `t.test(x, mu=mu0, alternative="two.sided"
/"less"/"greater")`

Exemplu:

O fabrică producătoare de baterii susține că durata medie de funcționare a produselor sale este de 180 ore. Pentru verificare, se aleg la întâmplare și se analizează 50 baterii. Duratele lor de funcționare se găsesc în dataframe-ul *Battery* din *PASWR*, în variabila *facilityA*. Vom testa la nivelul de semnificație 0.05 cele susținute de firmă.

Compararea mediilor a două populații

Pentru populații cu distribuție aproximativ normală sau volum de eșantion mare:

- ▶ testul t pentru eșantioane independente
- ▶ testul t pentru eșantioane dependente (perechi)

În cazul populațiilor cu distribuție non-normală și eșantioane mici se folosesc teste neparametrice:

- ▶ eșantioane independente → testul Mann - Whitney;
- ▶ eșantioane dependente (perechi) → testul Wilcoxon.

Testul t pentru două eșantioane independente

Se folosește pentru compararea mediilor a două populații independente cu distribuție aproximativ normală sau volume de eșantion suficient de mari.

- ▶ $H_0 : \mu_1 - \mu_2 = \mu_0$
- ▶ $H_a : \mu_1 - \mu_2 \neq \mu_0$ sau $H_a : \mu_1 - \mu_2 < \mu_0$ sau $H_a : \mu_1 - \mu_2 > \mu_0$

- ▶ Statistica test:
$$T = \frac{(\bar{X}_1 - \bar{X}_2) - \mu_0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

În ipoteza că H_0 e adevărată, T are distribuție t (Welch)

- ▶
$$t_{obs} = \frac{(\bar{X}_1 - \bar{X}_2) - \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$
- ▶ Dacă p-valoarea corespunzătoare lui t_{obs} este $< \alpha$, se respinge H_0 ; altfel, nu se respinge H_0 la nivelul de semnificație ales.

Testul t pentru două eșantioane independente

În R: `t.test(x, y, mu=mu0, alternative="two.sided"
/"less"/"greater")`

Exemplu:

Pentru setul de date *mtcars* din *datasets*, vom testa la nivelul de semnificație 0.05 ipoteza că mașinile cu transmisie manuală ($am = 0$) sunt mai puțin eficiente decât cele cu transmisie automată ($am = 1$) în privința consumului.

Testul t pentru două eșantioane dependente (perechi)

Se folosește pentru compararea mediilor a două populații cu distribuție aproximativ normală sau volume de eșantion suficient de mari, în cazul în care acestea constituie observații repetate sau perechi.

- ▶ $H_0 : \mu_1 - \mu_2 = \mu_0$
- ▶ $H_a : \mu_1 - \mu_2 \neq \mu_0$ sau $H_a : \mu_1 - \mu_2 < \mu_0$ sau $H_a : \mu_1 - \mu_2 > \mu_0$
- ▶ Se consideră $D_i = X_i - Y_i$, $i = \overline{1, n}$. Atunci D_i sunt i.i.d. și

$$\mu_1 - \mu_2 = \mu_0 \Leftrightarrow \mu_D = \mu_0.$$

→ se poate utiliza testul t pentru un singur eșantion reprezentând mulțimea diferențelor din perechile de observații.

Testul t pentru două eșantioane dependente (perechi)

În R: `t.test(x, y, mu=mu0, alternative="two.sided"
/"less"/"greater", paired=TRUE)`

Exemplu:

O dietă este promovată ca fiind capabilă să reducă în mod considerabil nivelul de glucoză din sânge. Zece pacienți diabetici sunt aleși aleator și puși să urmeze dieta o lună, și rezultatele lor sunt prezentate mai jos.

| | | | | | | | | | | |
|---------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| înainte | 268 | 225 | 252 | 192 | 307 | 228 | 246 | 298 | 231 | 185 |
| după | 106 | 236 | 253 | 110 | 203 | 101 | 211 | 176 | 194 | 203 |

Verificăm la nivelul de semnificație 0.05 dacă există suficiente dovezi care să susțină eficiența dietei la pacienții diabetici.

Testele Wilcoxon și Mann-Whitney

Date două populații X și Y , testele Wilcoxon și Mann-Whitney se folosesc pentru a verifica ipoteza

$$H_0 : P(X < Y) = \frac{1}{2} \text{ (valorile din cele două grupuri sunt similare)}$$

cu una dintre alternativele

- ▶ $P(X < Y) \neq \frac{1}{2}$
- ▶ $P(X < Y) < \frac{1}{2}$
- ▶ $P(X < Y) > \frac{1}{2}$.

Testele Wilcoxon și Mann-Whitney se pot folosi și pentru date ordinale.

În R:

```
wilcox.test(x, y, alternative="two.sided"/"less"/  
"greater", paired=T/F)
```

Testele Wilcoxon și Mann-Whitney

Exemplu:

Se planifică un studiu pilot care să testeze eficiența administrării de suplimente de vitamina E pentru prevenirea bolii Alzheimer. 20 subiecți cu vârste peste 65 ani sunt repartizați aleator în două grupuri. Primul grup (10 persoane) primește 400 UI/zi vitamina E, iar grupul al doilea primește un tratament placebo. Se înregistrează nivelul inițial de vitamina E în fiecare grup, obținându-se valorile:

Grup 1 : 7.5, 12.6, 3.8, 20.2, 6.8, 403.3, 2.9, 7.2, 10.5, 205.4

Grup 2 : 8.2, 13.3, 102.0, 12.7, 6.3, 4.8, 19.5, 8.3, 407.1, 10.2

Analizăm dacă există diferențe între grupuri la momentul inițial.