

## Laborator 6 - Statistică inferențială

### I. Inferență asupra dispersiilor a două populații - Testul F (simetric) asupra dispersiilor a două populații

Se consideră o două populații statistice cărora nu li se cunosc dispersiile  $\sigma_1^2$  și  $\sigma_2^2$ . Din cele două populații se extrag două eșantioane aleatoare simple (și independente între ele) cărora li se calculează dispersiile  $s_1^2$  și  $s_2^2$ . Scorul  $F = \frac{s_1^2}{s_2^2}$  este distribuit  $F(n_1 - 1, n_2 - 1)$ .

Testul  $F$  decurge astfel:

1. se formulează ipoteza nulă, care susține că dispersiile celor două populații sunt egale:

$$H_0 : \frac{\sigma_1^2}{\sigma_2^2} = 1$$

2. se formulează ipoteza alternativă care susține că dispersiile sunt diferite:

$$H_a : \frac{\sigma_1^2}{\sigma_2^2} \neq 1$$

3. se fixează nivelul de semnificație:  $\alpha$  (care uzual poate fi 1% sau 5%);

4. se calculează scorul testului:

$$F = \frac{s_1^2}{s_2^2}$$

5. se determină valorile critice  $F_s^*$  și  $F_d^*$ :

$$F_s^* = qf(\alpha/2, n_1 - 1, n_2 - 1)$$

$$F_d^* = qf(1 - \alpha/2, n_2 - 1, n_1 - 1) = 1/F_s^*$$

6. ipoteza nulă  $H_0$  este respinsă și se acceptă că dispersiile sunt diferite dacă

$$F \notin (F_s^*, F_d^*)$$

altfel nu există suficiente dovezi pentru a accepta că dispersiile sunt diferite.

**Exercițiu rezolvat.** Rezultatele unui test psihologic efectuat pe două eșantioane, unul de femei și unul de bărbați sunt următoarele:

bărbați:  $n_1 = 120, s_1 = 5.05$

femei:  $n_2 = 135, s_2 = 5.44$

Se poate trage concluzia că dispersiile celor două populații diferă semnificativ (1%)?

```
> alfa = 0.01
> n1 = 120
> n2 = 135
> s1 = 5.05
> s2 = 5.44
> critical_F_s = qf(alfa/2, n1 - 1, n2 - 1)
> critical_F_d = qf(1 - alfa/2, n2 - 1, n1 - 1)
> critical_F_s
> critical_F_d
> F_score
```

Rezultatul va fi  $F = 0.86175$ ,  $F_s^* = 0.62843$  și  $F_d^* = 1.59125$ ; deoarece  $F \in [F_s^*, F_d^*]$  ipoteza nulă nu poate fi respinsă. În acest caz putem considera că nu există dovezi semnificative pentru a afirma că dispersiile sunt diferite.

### Exerciții propuse

I.1 Scrieți o funcție (numită **F\_test**) care să calculeze și să returneze valorile critice și scorul testului  $F$  (parametrii funcției vor fi: tipul ipotezei alternative,  $\alpha$ ,  $n_1$ ,  $n_2$ ,  $\sigma_1$  etc., eșantioanele se pot extrage din fișier ca mai jos). Funcția aceasta va fi utilizată apoi la rezolvarea exercițiilor care urmează.

```
> x1 = read.table("program.txt", header = TRUE)[['A']]
> x2 = read.table("program.txt", header = TRUE)[['B']]
> n1 = length(x1)
> s1 = sd(x1)
> ...
```

I.2 Un profesor crede că un anumit program de lectură îmbunătățește abilitățile și dorința copiilor de a citi. Pentru aceasta el alege două grupuri de elevi: unul de 22 de elevi care urmează programul prescris ( $A$ ) și unul de 24 de elevi care nu urmează acest program ( $B$ ). Rezultatele sunt date în fișierul *program.txt*.

Să se decidă cu 1% și 5% nivel de semnificație dacă dispersiile celor două populații sunt diferite.

I.3 Cercetătorii studiază amplitudinea mișcării obținută prin stimularea nervoasă a șoarecilor. Pentru șoarecii drogați se obțin următoarele date:

12.512 12.869 19.098 15.350 13.297 15.589

Pentru șoarecii normali se obțin următoarele date:

11.074 9.686 12.164 8.351 12.182 11.489

Influența drogurilor este semnificativă în ceea ce privește cele două dispersii (5% nivel de semnificație)?

## II. Inferență asupra mediilor a două populații - Testul $t$ pentru diferența mediilor unor populații cu dispersii necunoscute

Se consideră o două populații statistice cărora nu li se cunosc dispersiile. Se aleg două eșantioane aleatoare simple și independente între ele cu mediile de selecție  $\bar{x}_{n_1}$  și  $\bar{x}_{n_2}$  și dispersiile  $\sigma_1^2$  și  $\sigma_2^2$ . Înaintea efectuării testului  $t$  se folosește testul  $F$  pentru a decide dacă dispersiile celor două populații sunt diferite. Dacă dispersiile nu sunt diferite

$$t = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$$

urmează o

distribuție Student<sup>1</sup>:  $t(n_1 + n_2 - 2)$ . Dacă dispersiile sunt diferite, atunci

$$t = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

urmează o distribuție Student:  $t(\min(n_1 - 1, n_2 - 1))$ .

Testul  $t$  decurge astfel:

---

<sup>1</sup>În acest caz  $s = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$

1. se formulează ipoteza nulă, care susține că diferența mediilor celor două populații ia o valoare particulară (de cele mai multe ori zero):

$$H_0 : \mu_1 - \mu_2 = m_0$$

2. se formulează o ipoteză alternativă care poate fi de trei feluri:

$$H_a : \mu_1 - \mu_2 < m_0 \quad (\text{ipoteză asimetrică la stânga}) \text{ sau}$$

$$H_a : \mu_1 - \mu_2 > m_0 \quad (\text{ipoteză asimetrică la dreapta}) \text{ sau}$$

$$H_a : \mu_1 - \mu_2 \neq m_0 \quad (\text{ipoteză simetrică})$$

3. se fixează nivelul de semnificație:  $\alpha$  (care uzual poate fi 1% sau 5%);  
4. se calculează scorul testului:

- a) dacă dispersiile sunt diferite:

$$t = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

iar numărul de grade de libertate este  $df = \min(n_1 - 1, n_2 - 1)$ .

- b) dacă dispersiile sunt "egale"

$$t = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$$

unde  $s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$ , iar numărul de grade de libertate este  $df = n_1 + n_2 - 2$ .

5. se determină valoarea critică  $t^*$ :

$$t^* = qt(2\alpha, df) \quad \text{pentru ipoteză } H_a \text{ asimetrică la stânga } (t^* < 0),$$

$$t^* = qt(1 - \alpha, df) \quad \text{pentru ipoteză } H_a \text{ asimetrică la dreapta } (t^* > 0),$$

$$t^* = -qt(\alpha/2, df) = qt(1 - \alpha/2, df) \quad \text{pentru ipoteză } H_a \text{ simetrică } (t^* > 0).$$

6. ipoteza nulă  $H_0$  este respinsă dacă

$$t < t^* \quad \text{pentru ipoteză } H_a \text{ asimetrică la stânga sau}$$

$$t > t^* \quad \text{pentru ipoteză } H_a \text{ asimetrică la dreapta sau}$$

$$|t| > |t^*| \quad \text{pentru ipoteză } H_a \text{ imetrică,}$$

altfel se spune că **nu există suficiente dovezi pentru a respinge ipoteza nulă  $H_0$  și a accepta ipoteza alternativă  $H_a$ .**



**Exercițiu rezolvat.** Rezultatele unui test psihologic efectuat pe două eșantioane, unul de femei și unul de bărbați sunt următoarele:

bărbați:  $n_1 = 110, \bar{x}_1 = 25.84, s_1 = 4.25$

femei:  $n_2 = 105, \bar{x}_2 = 21.53, s_2 = 3.85$

Se poate trage concluzia că mediile celor două populații diferă semnificativ (1%)?

*Observație.* Testarea diferenței mediilor echivalează cu o ipoteză alternativă de tipul  $\mu_1 - \mu_2 \neq 0 = m_0$ .

```
> alfa = 0.01
> m0 = 0
> sample1_mean = 25.84
> sample2_mean = 21.53
> n1 = 110
> n2 = 105
> s1 = 4.25
> s2 = 3.85
> critical_F_s = qf(alfa/2, n1 - 1, n2 - 1) # testul F
> critical_F_d = qf(1 - alfa/2, n2 - 1, n1 - 1)
> F_score = s1^2 / s2^2
> if(F_score < critical_F_s | F_score > critical_F_d) {
+ df = min(n1 - 1, n2 - 1)
+ combined_s = sqrt(s1^2/n1 + s2^2/n2)
+ } else {
+ df = n1 + n2 - 2
+ combined_s = sqrt(((n1 - 1)*s1^2 + (n2 - 1)*s2^2)/df)*sqrt(1/n1+1/n2)
+ }
> critical_t = qt(1 - alfa/2, df)
> t_score = (sample1_mean - sample2_mean - m0)/combined_sigma
> critical_t
> t_score
```

Rezultatul va fi  $|t^*| = 2.59910 > |t| = 1.87248$ , ipoteza nulă nu poate fi respinsă.

### Exerciții propuse

II.1 Scrieți o funcție (numită, de exemplu, **t\_test\_means**) care să calculeze și să returneze valoarea critică și scorul testului  $t$  pentru diferența mediilor; această funcție va include și testul F simetric (parametrii vor fi  $\alpha, n_1, n_2, s_1$  etc.). Funcția aceasta va fi utilizată, pentru rezolvarea exercițiilor care urmează.

I.2 Un profesor crede că un anumit program de lectură îmbunătățește abilitățile și dorința copiilor de a citi. Pentru aceasta el alege două grupuri de elevi: unul de 22 de elevi care urmează programul prescris (A) și unul de 24 de elevi care nu urmează acest program (B). Rezultatele sunt date în fișierul *program.txt*.

Să se decidă cu 1% și 5% nivel de semnificație dacă mediile celor două populații diferă semnificativ.

I.3 Cercetătorii studiază amplitudinea mișcării obținută prin stimularea nervoasă a șoarecilor. Pentru șoarecii drogați se obțin următoarele date:

12.512 12.869 19.098 15.350 13.297 15.589

Pentru șoarecii normali se obțin următoarele date:

11.074 9.686 12.164 8.351 12.182 11.489

Influența drogurilor este semnificativă în ceea ce privește cele două medii (1% nivel de semnificație)?

### III. Inferență asupra mediei - Testul Z pentru diferența mediilor unor populații cu dispersii cunoscute

Se consideră o două populații statistice cărora li se cunoaște dispersiile  $\sigma_1^2$  și  $\sigma_2^2$ . Se aleg două eșantioane aleatoare simple și independente între ele cu mediile de selecție  $\bar{x}_{n_1}$  și  $\bar{x}_{n_2}$ . Dacă populațiile urmează o lege normală sau dimensiunea eșantioanelor este suficient de mare, scorul

$$z = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

este distribuit (eventual cu aproximație) normal standard:  $N(0, 1)$ .

Testul Z decurge astfel:

1. se formulează ipoteza nulă, care susține că diferența mediilor celor două populații ia o valoare particulară:

$$H_0 : \mu_1 - \mu_2 = m_0$$

2. se formulează o ipoteză alternativă care poate fi de trei feluri:

$$H_a : \mu_1 - \mu_2 < m_0 \quad (\text{ipoteză asimetrică la stânga}) \text{ sau}$$

$$H_a : \mu_1 - \mu_2 > m_0 \quad (\text{ipoteză asimetrică la dreapta}) \text{ sau}$$

$$H_a : \mu_1 - \mu_2 \neq m_0 \quad (\text{ipoteză simetrică})$$

3. se fixează nivelul de semnificație:  $\alpha$  (care uzual poate fi 1% sau 5%);
4. se calculează scorul testului:

$$\frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

5. se determină valoarea critică  $z^*$ :

$$z^* = qnorm(\alpha, 0, 1) \quad \text{pentru ipoteză } H_a \text{ asimetrică la stânga,}$$

$$z^* = qnorm(1 - \alpha, 0, 1) \quad \text{pentru ipoteză } H_a \text{ asimetrică la dreapta,}$$

$$z^* = qnorm(1 - \alpha/2, 0, 1) \quad \text{pentru ipoteză } H_a \text{ simetrică.}$$

6. ipoteza nulă  $H_0$  este respinsă dacă

$$z < z^* \quad \text{pentru ipoteză } H_a \text{ asimetrică la stânga sau}$$

$$z > z^* \quad \text{pentru ipoteză } H_a \text{ asimetrică la dreapta sau}$$

$$|z| > |z^*| \quad \text{pentru ipoteză } H_a \text{ simetrică,}$$

altfel se spune că **nu există suficiente dovezi pentru a respinge ipoteza nulă  $H_0$  și a accepta ipoteza alternativă  $H_a$ .**

**Exercițiu rezolvat.** Se compara durata de viața a doua tipuri de baterii. Primul tip are o deviație standard de 4 ore, al doilea tip are o deviație standard de 3 ore. Se aleg două eșantioane fiecare de dimensiune de 100 de baterii. Pentru primul eșantion media de viața este de 48 de ore, iar pentru cel de-al doilea de 47 de ore.

Să se testeze diferența mediilor de viața cu 5% nivel de semnificație.

*Observație.* Testarea diferenței mediilor echivalează cu o ipoteză alternativă de tipul  $\mu_1 - \mu_2 \neq 0 = m_0$ .

```
> alfa = 0.05
> m0 = 0
> sample1_mean = 48
> sample2_mean = 47
> n1 = 100
> n2 = 100
> sigma1 = 4
> sigma2 = 3
> combined_sigma = sqrt(sigma1^2/n1 + sigma2^2/n2)
> critical_z = qnorm(1 - alfa/2)
> z_score = (sample1_mean - sample2_mean - m0)/combined_sigma
> critical_z
> z_score
```

Rezultatul va fi  $|z^*| = 1.95996 < |z| = 2.00$ , ipoteza nulă va fi respinsă și se acceptă că mediile celor două populații sunt diferite.

### Exerciții propuse

III.1 Scrieți o funcție (numită, de exemplu, **z\_test\_means**) care să calculeze și să returneze valoarea critică și scorul testului  $z$  pentru diferența mediilor (parametrii vor fi  $\alpha$ ,  $n_1$ ,  $n_2$ ,  $\sigma_1$  etc.). Funcția aceasta va fi utilizată, pentru rezolvarea exercițiilor de mai jos.

III.2 80 dintre angajații unei firme au un salariu mediu săptămânal de 160\$ (deviația standard a întregii populații fiind 3.24\$). 70 dintre angajații unei alte firme au în medie 155\$ salariu pe săptămână (deviația standard a întregii populații fiind 2.25\$). Să se testeze dacă salariul mediu săptămânal la cele două firme diferă semnificativ (1% nivel de semnificație).

III.3 Un raport recent arată că absolvenții de universitate fără diplomă se căsătoresc mai repede decât cei cu diplomă. Sunt aleși câte 100 de indivizi din cele două populații; pentru aceste două eșantioane absolvenții fără diplomă se căsătoresc în medie la 22.8 ani (deviația standard cunoscută populației fiind  $\sigma_1 = 1.3$  ani) iar cei cu diplomă la 23.3 ani (deviația cunoscută a populației este  $\sigma_2 = 1.9$  ani).

Cu 1% nivel de semnificație se poate trage concluzia că raportul este corect?

III.4 O firmă producătoare de pastă de dinți dorește să știe dacă introducerea unui nou compus chimic mărește performanțele produsului său. Sunt alese două eșantioane fiecare a câte 100 de indivizi. Pentru primul grup (care folosește noul produs) se numără în medie 3 carii (deviația standard cunoscută este 0.6), pentru cel de-al doilea grup se numără în medie 3.5 carii (deviația standard este de 0.4).

Cu 1% nivel de semnificație se poate trage concluzia că prin introducerea noului compus pasta de dinți își îmbunătățește performanțele? Dar cu 5%?

III.5 Pentru două eșantioane provenind din două populații statistice distincte  $A$  și  $B$  se determină următoarele valori

$$(A) \quad \sigma_1 = 0.75, n_1 = 155, \bar{x}_{n_1} = 15$$



$$(B) \sigma_2 = 0.78, n_2 = 150, \bar{x}_{n_2} = 14.5$$

Să se testeze diferența mediilor celor două populații cu 1% și cu 5% nivel de semnificație.