

Laborator 7 - Statistică inferențială

I. Testarea ipotezelor statistice - Testul χ^2 pentru dispersia unei populații

Se consideră o populație statistică normal distribuită căreia nu i se cunoaște dispersia σ^2 . Pentru un eșantion aleator simplu cu dimensiunea n și dispersia s^2 , scorul $\chi^2 = (n - 1) \frac{s^2}{\sigma^2}$ este distribuit $\chi^2(n - 1)$. Testul decurge astfel:

1. se formulează ipoteza nulă, care susține că dispersia populației ia o valoare particulară:

$$H_0 : \sigma^2 = \sigma_0^2$$

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

$$H_a : \sigma^2 < \sigma_0^2 \quad (\text{ipoteză asimetrică la stânga}) \text{ sau}$$

$$H_a : \sigma^2 > \sigma_0^2 \quad (\text{ipoteză asimetrică la dreapta}) \text{ sau}$$

$$H_a : \sigma^2 \neq \sigma_0^2 \quad (\text{ipoteză simetrică})$$

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

4. se calculează scorul testului:

$$\chi^2 = (n - 1) \cdot \frac{s^2}{\sigma_0^2}$$

5. se determină valorile critice:

$$\chi^{2*} = qchisq(\alpha, n - 1) \quad \text{pentru ipoteză } H_a \text{ asimetrică la stânga,}$$

$$\chi^{2*} = qchisq(1 - \alpha, n - 1) \quad \text{pentru ipoteză } H_a \text{ asimetrică la dreapta,}$$

$$\chi_s^{2*} = qchisq\left(\frac{\alpha}{2}, n - 1\right), \chi_d^{2*} = qchisq\left(1 - \frac{\alpha}{2}, n - 1\right) \quad \text{pentru } H_a \text{ simetrică.}$$

6. ipoteza nulă H_0 este respinsă dacă

$$\chi^2 < \chi^{2*} \quad \text{pentru ipoteză } H_a \text{ asimetrică la stânga sau}$$

$$\chi^2 > \chi^{2*} \quad \text{pentru ipoteză } H_a \text{ asimetrică la dreapta sau}$$

$$\chi^2 \notin (\chi_s^{2*}, \chi_d^{2*}) \quad \text{pentru ipoteză } H_a \text{ simetrică,}$$

dacă nu suntem într-una din aceste situații, atunci 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 de inteligență efectuat pe un eșantion aleator simplu de 120 de indivizi dintr-o populație distribuită normal oferă o deviație standard de 13 puncte. Se poate trage concluzia că dispersia populației este mai mică decât 225? (1%)

```

> alfa = 0.01
> n = 120
> s_square = 169
> sigma_square = 225
> critical_Chi_square = qchisq(alfa, s - 1)
> Chi_square_score = (n - 1)*s_square / sigma_square
> critical_Chi_square
> Chi_square_score

```

Rezultatul va fi $\chi^2 = 89.3822 > \chi^{2*} = 86.07383$, deci ipoteza nulă nu poate fi respinsă.

Exerciții propuse

- I.1 Scrieți o funcție (numită **ChiSq_Variance_test**) care să calculeze și să returneze valoarea critică și scorul testului χ^2 pentru dispersie (parametrii funcției vor fi: α , n , s^2 , σ_0^2). Funcția aceasta va fi utilizată apoi pentru rezolvarea exercițiilor care urmează.
- I.2 O companie de cosmetice afirmă că dispersia cantității de parfum într-un flacon este aproximativ egală cu 0.2. Se consideră un eșantion de dimensiune 100, pentru care se calculează o dispersie $s^2 = 0.22$. Să se testeze cu 5% nivel de semnificație dacă dispersia este mai mare decât 0.2.
- I.3 Pentru un eșantion de dimensiune 35, se măsoară o deviație standard de 10.5. Să se testeze cu nivel de semnificație de 1% dacă deviația standard a populației este diferită de 11.
- I.4 Să se testeze ipoteza că dispersia unei populații este diferită de 12.8, dacă pentru un eșantion dat de dimensiune 49, dispersia măsurată este mai mică decât 13.
- I.5 Pentru 36 de pachete de țigarete este găsită o deviație standard a concentrației de nicotină $s = 0.3$ mg. Să se testeze cu 1% nivel de semnificație dacă dispersia este mai mare decât 0.0625.

II. Testarea ipotezelor statistice - Testul χ^2 pentru experimente multinomiale

Se consideră un experiment cu s rezultate posibile, cu probabilități p_1, p_2, \dots, p_s . Experimentul se repetă în mod independent de n ori și se notează cu O_i numărul de realizări ale rezultatului i (O_i este o variabilă aleatoare distribuită binomial $B(n, p_i)$):

$$\sum_{i=1}^s O_i = n.$$

Pentru rezultatul i numărul mediu de realizări este (media variabilei O_i) $E_i = np_i$.

Statistica $\sum_{i=1}^s \frac{(O_i - E_i)^2}{E_i}$ urmează o distribuție $\chi^2(s - 1)$. Testul multinomial inferează asupra probabilităților p_1, p_2, \dots, p_s și decurge astfel:

1. se formulează ipoteza nulă, care susține că probabilitățile iau valori fixate.

$$H_0 : p_i = \pi_i, \forall i$$

2. se formulează o ipoteză alternativă:

$$H_a : p_i \neq \pi_i, \text{ pentru măcar un } i$$

3. se fixează nivelul de semnificație: α (1% sau 5%);

4. se calculează scorul testului:

$$\chi^2 = \sum_{i=1}^s \frac{(O_i - E_i)^2}{E_i}$$

5. se calculează valoarea critică:

$$\chi^{2*} = qchisq(1 - \alpha, s - 1)$$

6. ipoteza nulă H_0 este respinsă dacă $\chi^2 > \chi^{2*}$, altfel se spune că **nu există suficiente dovezi pentru a respinge ipoteza nulă H_0 și a accepta ipoteza alternativă H_a .**

Exercițiu rezolvat. Studenții doresc o cât mai mare libertate în alegerea cursurilor. Șapte cursuri similare, predate de profesori diferiți, au fost alese de 119 studenți astfel (ordinea este aleatoare):

18 12 25 23 8 19 14

Cu 5% nivel de semnificație, se poate trage concluzia că studenții au avut preferințe în alegerea cursurilor (i. e., există măcar un curs care este ales cu probabilitate diferită de 1/7)?

```
> alfa = 0.01
> O = c(18, 12, 25, 23, 8, 19, 14)
> p = c(1/7, 1/7, 1/7, 1/7, 1/7, 1/7, 1/7)
> s = length(O)
> n = sum(O)
> E = n*p
> critical_Chi_square = qchisq(1 - alfa, s - 1)
> Chi_square_score = sum((O - E)^2/E)
> critical_Chi_square
> Chi_square_score
```

Rezultatul va fi $\chi^2 = 12.94118 > \chi^{2*} = 12.59159$, deci ipoteza nulă poate fi respinsă, măcar unul dintre cursuri este preferat față de un altul.

Exerciții propuse

II.1 Scrieți o funcție (numită **multinomial_test**) care să calculeze și să returneze valoarea critică și scorul testului multinomial (parametrii funcției vor fi: α , O și p). Funcția aceasta va fi utilizată apoi pentru rezolvarea exercițiilor care urmează.

II.2 Compania Mars Inc. pretinde că bomboanele M&M sunt distribuite pe culori astfel: 35% sunt maro, 20% galbene, 15% roșii, 10% portocalii, 10% verzi și 10% albastre. Pentru 100 de bomboane, se observă următoarea distribuție pe culori:

	Maro	Galben	Roșu	Portocaliu	Verde	Albastru
Frecvența observată	33	26	18	9	7	7

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

II.3 Dintre șoferii ce au suferit un accident de circulație, au fost selectați 88 de șoferi, în mod aleatoriu, și au fost împărțiți în funcție de vârstă.

Vârstă	sub 25 de ani	între 25 și 40	între 40 și 65	peste 65 de ani
Frecvența observată	21	36	12	19

Testați ipoteza că distribuția accidentelor pe categorii de vârstă este următoarea: 16%, 44%, 27%, 13%.