

# PRELUCRAREA DATELOR DE SONDAJ SI INFERENȚA STATISTICĂ

## Testarea ipotezelor statistice

Stud. Master - AMP

ISAIC- MANIU ALEXANDRU  
web [www.amaniu.ase.ro](http://www.amaniu.ase.ro)  
e-mail AL.ISAIC-MANIU@CSIE.ASE.RO

17.XI.2013

1

## Cateva elemente recapitulative

- **Estimare**- operație de stabilire, în baza datelor unui eșantion, a valorilor parametrilor repartiției populației din care a fost prelevat eșantionul
- ✓ Rezultatul, se poate exprima printr-o valoare unică (**estimator punctual**), sau printr-un interval( numit frecvent “**interval de incredere** ”)
- ✓ Utilizand doar părți din populație, rezultatele obtinute sunt acompaniate de anumite **riscuri**
- ❖ Spre deosebire de **statistica descriptivă**, inferența folosește procedee specifice bazate pe modele matematice (în esență, probabiliste) pentru analiza materialului statistic organizat de metodele descriptive

## Medie aritmetică de sondaj -

(sample average, sample mean-value)

- Raportul dintre suma tuturor valorilor x observate în eșantionul considerat și numărul total n al acestora:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i^*$$

- Observații  
În cazul valorilor observate, aranjate în ordine crescătoare sau descrescătoare:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n n_i x_i = \sum_{i=1}^n f_i x_i$$

- în care:  
n numărul total al valorilor observate;  
n<sub>i</sub> frecvența absolută corespunzătoare valorii x<sub>i</sub>  
f<sub>i</sub> frecvența relativă corespunzătoare valorii x<sub>i</sub>

17.XI.2013

3

## Media caracteristicii binare (alternative)

Când valorile variabilei aleatoare  $X$  sunt proporția elementelor  $A$  sau respectiv a elementelor  $\text{non } A$  (de exemplu, proporția cetățenilor cu intenție de a se prezenta la vot și absenteiștii), atunci valorile tipice respective sunt:

media:  $M(X) = p$

dispersia de sondaj:  $\text{Var}^2(X) = pqn$

În care:

$n$  – volumul eșantionului;

$p$  și  $q$  – proporțiile respective ale elementelor  $A$  și  $\text{non } A$ .

17.XI.2013

4

## Dispersia de sondaj - ( $s^2$ )

- Momentul centrat de ordinul doi:

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

- Valoarea numerică a acestui indicator sintetic caracterizează împrăștierea repartiției statistice
- Dispersia de sondaj poate fi folosită ca estimare aproximativă a dispersiei din populația originală, considerându-se formula corectă:

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

- **Eroarea în estimarea** unui parametru ( $\hat{\theta} - \theta$ ) unde  $\hat{\theta}$  este rezultatul estimării, iar  $\theta$  este parametrul estimat

17.XI.2013

5

## Interval de încredere

(Confidence interval)

- **1. Interval de încredere unilateral** (*One-sided confidence interval; Intervalle de confiance unilatéral*). Dacă  $Z$  este o funcție a valorilor observate, iar  $\theta$  este un parametru de estimat al populației, astfel ca probabilitatea este cel puțin egală cu o probabilitate fixată de cercetător
- intervalul cuprins între cea mai mică valoare posibilă a lui  $\theta$  și  $Z$  (sau intervalul între  $Z$  și cea mai mare valoare posibilă a lui  $\theta$ ) este **intervalul unilateral de probabilitate  $P$  pentru  $\theta$** .

17.XI.2013

6

## 2. Interval de încredere bilateral

- Dacă  $Z_1$  și  $Z_2$  sunt două funcții ale valorilor observate, iar  $\theta$  este un parametru estimat al populației, astfel ca probabilitatea este cel puțin egală cu  $1 - \alpha$ , [unde  $1 - \alpha$  este un număr fixat, pozitiv și mai mic decât 1], intervalul dintre  $Z_1$  și  $Z_2$  este un interval de încredere bilateral de pentru  $\theta$
- Limitele  $Z_1$  și  $Z_2$  ale intervalului de încredere sunt statistici care, în general, au valori diferite de la un eșantion la altul

17.XI.2013

7

---

---

---

---

---

---

---

---

---

---

## Erori in verificarea ipotezelor statistice

- ❖ Eroare de genul întâi : ipoteza H se respinge, când ea este adevărată
- ❖ Eroare de genul al doilea: ipoteza H se admite, când ea este falsă

Probabilitățile de a fi comise erori sunt:

probabilitatea erorii de genul întâi

✓ risc de genul I ( $\alpha$ ) și respectiv

probabilitatea erorii de genul al doilea-

✓ risc de genul II ( $\beta$ )

17.XI.2013

8

---

---

---

---

---

---

---

---

---

---

## Ipoteza nulă și ipoteza alternativă

- Afirmații asupra unuia sau mai multor parametri, sau asupra unor repartiții, care urmează a fi validate prin teste statistice.
- Decizia asupra ipotezei nule este luată pe baza unui test statistic.
- Testul este construit cu elemente aleatoare, iar decizia comportă un anumit risc de eroare.

$$H_0 : (p_1 = p_2) \quad H_1 : (p_1 \neq p_2)$$

- Ipoteza nulă ( $H_0$ ) se referă la afirmații supuse testării, în timp ce ipoteza alternativă ( $H_1$ ) se referă la afirmații care vor fi acceptate dacă se respinge ipoteza nulă.

17.XI.2013

9

---

---

---

---

---

---

---

---

---

---

### ✦ Test statistic (Statistical test)

- ◊ Procedura statistică prin care se decide dacă ipoteza nulă poate fi respinsă în favoarea ipotezei alternative sau nu.
- ◊ În general, un test preia apriori o anumită ipoteză, care trebuie verificată (de exemplu, ipoteza de independență a observațiilor, ipoteza de normalitate, ipoteza egalității unor medii etc.).

### Test "neparametric" (Distribution-free test)

Testul în care funcția de repartiție a statisticii decizionale utilizate **nu depinde de funcția de repartiție a observațiilor**. Sensul este cel dat de termenul englezesc. Denumirea *neparametric* a fost aleasă mai curând pentru ușurința exprimării. În română ar trebui să spunem "*test independent de repartiția inițială a ...*", dar fiind o formulare prea lungă s-a optat pentru neparametric.

17.XI.2013

10

## TESTE PENTRU EGALITATEA MEDIILOR

### Testul U

- Test utilizat pentru verificarea ipotezelor referitoare la mediile populațiilor normale când se cunosc dispersiile teoretice.
- Testul U are forme diferite, în funcție de ipotezele statistice formulate:
- Ex. :a) se verifică ipoteza  $H_0: m = m_0$ , testul U are expresia:

$$U = \frac{\bar{x} - m_0}{\frac{\sigma}{\sqrt{n}}}$$

17.XI.2013

11

### ■ Test U

- ◊ Test utilizat pentru verificarea ipotezelor referitoare la mediile populațiilor normale când se cunosc dispersiile teoretice.
- b) Când se verifică ipoteza egalității a două medii corespunzând la două populații normale care au aceeași dispersie teoretică  $\sigma^2$ , testul U are expresia:

$$U = \frac{\bar{x}_1 - \bar{x}_2}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

17.XI.2013

12

c) Când se verifică ipoteza egalității a două medii, corespunzând la două populații normale care au dispersiile teoretice cunoscute, însă negale, testul U are expresia:

$$U = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

În expresiile de mai sus,  $\bar{x}_1$  și  $\bar{x}_2$  sunt mediile aritmetice de sondaj, iar  $n_1$  și  $n_2$  sunt volumele eșantioanelor prelevate din populația originală.

4 Regiunea critică a testului U este  $|U| > u(\alpha)$ , în care  $u(\alpha)$  se citește din tabelul repartiției normale normale  $N(0,1)$ , astfel încât:

$$P[|U| > u(\alpha)] = 1 - \frac{1}{\sqrt{2\pi}} \int_{-u(\alpha)}^{u(\alpha)} e^{-\frac{z^2}{2}} dz = \alpha$$

17.XI.2013

13

## Testul HI - patrat

- Regiunea critică a testului  $\chi^2$  pentru verificarea ipotezei  $p_1 = p_2 = \dots = p_m$  se construiește pe baza indicatorului statistic de forma:

$$\chi^2 = \sum_{i=1}^n \frac{(n_i - np_i)^2}{np_i}$$

- care pentru  $n \rightarrow \infty$  are repartiția  $\chi^2$  cu  $k - 1$  grade de libertate

17.XI.2013

14

## Test F – Snedecor

- Testul statistic în care, pentru validarea ipotezei nule, statistica utilizată presupune existența repartiției F.
- Testul este utilizat pentru verificarea ipotezei egalității dispersiilor de sondaj obținute în două eșantioane independente

Statistica testului F este definit prin relația:

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

- astfel încât  $s_1^2 \geq s_2^2$  în care  $s_1^2$  și  $s_2^2$  sunt dispersiile de sondaj ale celor două eșantioane.

17.XI.2013

15

**Repartiție t** (t – Student's distribution)

Repartiția de probabilitate a unei variabile aleatoare continue, care are funcția de densitate de probabilitate exprimată prin:

$$f(t; v) = \frac{1}{\sqrt{\pi v}} \left( \frac{\Gamma\left(\frac{v+1}{2}\right)}{\Gamma\left(\frac{v}{2}\right)} \right) \left( \frac{1}{1 + \frac{t^2}{v}} \right)^{\frac{v+1}{2}}$$

unde  $-\infty < t < +\infty$  de parametru  $v = 1, 2, \dots$  iar  $\Gamma$  este funcția gamma.

**Observație**

Raportul dintre două variabile aleatoare independente, numărătorul fiind o variabilă normală normată și numitorul fiind rădăcina pătrată pozitivă a raportului dintre o variabilă aleatoare  $\chi^2$  și numărul său de grade de libertate  $v$ , este o repartiție Student cu  $v$  grade de libertate. Dacă variabila aleatoare  $X$  are o repartiție Student cu  $v > 2$  grade de libertate, atunci:

$$M(X) = 0; \text{Var}(X) = \frac{v}{v-2}$$

17.XI.2013

16

**Testul t – STUDENT**

◇ Testul statistic în care, pentru validarea ipotezei nule, statistica utilizată presupune existența repartiției t (Student).

◇ Testul este aplicat, de exemplu, la următoarele probleme:

a. când se verifică ipoteza  $H_0: m = m_0$ , indicatorul t are expresia:

$$t = \frac{\bar{x}_1 - m_0}{\frac{s}{\sqrt{n}}}$$

cu  $v = n - 1$  grade de libertate, n fiind volumul eșantionului.

b. când se verifică ipoteza egalității a două medii corespunzând la două populații normale care au aceeași dispersie teoretică (necunoscută), indicatorul t are expresia:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

17.XI.2013

cu  $v = n_1 + n_2 - 2$  grade de libertate.

17

c. când se verifică ipoteza egalității a două medii de sondaj corespunzând la două populații normale care au dispersiile teoretice neegale, indicatorul t are expresia:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

această expresie are un număr

de  $v$  grade de libertate, rezultat din formula:

$$v = \frac{1}{\frac{c^2}{n_1 - 1} + \frac{(1 - c^2)}{n_2 - 1}}; c = \frac{\frac{s_1^2}{n_1}}{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Regiunea critică a testului t, în toate cazurile menționate, este:

$$|t| > t_v(\alpha)$$

◇ Pentru  $v > 30$ , testul t poate fi înlocuit cu testul U.

17.XI.2013

18

Equaliteto o  
dva proporcije

$$Z_c = \frac{|p_1 - p_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

CANDIDAT - A

CURS =  $p_1 = 52\%$ ;  $n_1 = 1150$   
CSOP =  $p_2 = 45\%$ ;  $n_2 = 1500$

$$Z_c = \frac{|0,52 - 0,45|}{\sqrt{\frac{0,25}{1150} + \frac{0,25}{1500}}} = 3,89$$

$Z_c > Z_{\text{tbl.}} \Rightarrow$  Diferencija  
• Efektivnost u post 36,8%!

17.XI.2013

19

---

---

---

---

---

---

---

---