

Statistica matematică

► Statistica matematică este o ramură a matematicii aplicate, care se ocupă de colectarea, gruparea, analiza și interpretarea datelor referitoare la anumite fenomene în scopul obținerii unor previziuni;

- statistica descriptivă: metode de colectare, organizare, sintetizare, prezentare și descriere a datelor numerice (sau nenumerice) într-o formă convenabilă

- statistica inferențială: metode de interpretare a rezultatelor obținute prin metodele statisticii descriptive, utilizate apoi pentru luarea deciziilor.

► O *colectivitate* sau *populație statistică* \mathcal{C} este o mulțime de elemente care au anumite însușiri comune ce fac obiectul analizei statistice. Numărul elementelor populației se numește volumul populației.

Exemple de populații statistice: mulțimea persoanelor dintr-o anumită țară, localitate, zonă etc. într-un anumit an, multimea gospodăriilor din România la un moment dat, mulțimea consumatorilor unui produs, mulțimea societăților care produc un anumit produs, angajații unei societăți, studenții unei facultăți.

► *Eșantionul* \mathcal{E} reprezintă o submulțime a unei populații statistice $\mathcal{E} \subset \mathcal{C}$, constituită după criterii bine stabilite:

a) să fie aleatoare;

b) toate elementele colectivității să aibe aceeași șansă de a fi alese în eșantion;

c) eșantionul să fie reprezentativ (structura eșantionului să fie apropiată de structura populației);

d) volumul eșantionului să fie suficient de mare.

► *Unitatea statistică* (indivizii) este elementul, entitatea de sine stătătoare a unei populații statistice, care posedă o serie de trăsături caracteristice ce-i conferă apartenența la populația studiată.

De exemplu: *unitatea statistică simplă*: un salariat, un student, un agent economic, o trăsătură, o părere; *unitatea statistică complexă*: o grupă de studenți sau o echipă de salariați, o familie sau o gospodărie, o categorie de mărfuri.

► *Variabila statistică* sau *caracteristica* reprezintă o însușire, o proprietate măsurabilă a unei unități statistice, întâlnită la toate unitățile care aparțin aceleiași colectivități și care prezintă variabilitate de la o unitate statistică la alta. Caracteristica sau variabila statistică corespunde unei variabile aleatoare.

Exemple de caracteristici: vârsta, salariul, preferințele politice, prețul unui produs, calitatea unor servicii, nivelul de studii.

- caracteristici continue (greutatea, înălțimea, valoarea glicemiei, temperatura aerului)

- caracteristici discrete (numări elevi ai unei școli, numărul liceelor existente într-un oraș)

- caracteristici nominale (culoarea ochilor, ramura de activitate, religia)

- caracteristici nominale ordonate (starea de sănătate / calitatea unor servicii - precară, mai bună, bună, foarte bună)

- caracteristici dichotomiale (binare) (stagiul militar - satisfăcut/nesatisfăcut, starea civilă - căsătorit/necăsătorit, genul - masculin/feminin)

► *Datele statistice* reprezintă observațiile rezultate dintr-o cercetare statistică, sau ansamblul valorilor colectate în urma unei cercetări statistice.

De exemplu: un angajat al unei companii are o vechime de 6 ani în muncă. Angajatul reprezintă unitatea statistică, vechimea în muncă este caracteristica (variabila) cercetată, iar 6 este valoarea acestei caracteristici.

O *colectivitate* (populație) \mathcal{C} este cercetată din punctul de vedere al caracteristicii (variabilei statistice) X .

► Fie $\mathcal{E} \subset \mathcal{C}$ un eșantion. Se numesc *date de selecție* relative la caracteristica X datele statistice x_1, \dots, x_n obținute prin cercetarea indivizilor care fac parte din eșantionul \mathcal{E} .

► Datele de selecție x_1, \dots, x_n pot fi considerate ca fiind valorile unor variabile aleatoare X_1, \dots, X_n , numite variabile de selecție și care se consideră a fi variabile aleatoare independente și având aceeași distribuție ca X .

► Fie x_1, \dots, x_n datele statistice pentru caracteristica cercetată X , notăm cu X_1, \dots, X_n variabilele de selecție corespunzătoare. Fie $g: \mathbb{R}^n \rightarrow \mathbb{R}$ o funcție astfel încât $g(X_1, \dots, X_n)$ este o variabilă aleatoare.

$g(X_1, \dots, X_n)$ se numește *funcție de selecție* sau *estimator*

$g(x_1, \dots, x_n)$ se numește valoarea funcției de selecție sau *valoarea estimatorului*.

Distribuția caracteristicii X poate fi

1) complet specificată (de ex.: $X \sim \text{Exp}(3)$, $X \sim \text{Bin}(10, 0.3)$, $X \sim N(0, 1)$)

2) specificată, dar depinzând de unul sau mai mulți parametri necunoscuți
(de ex.: $X \sim \text{Exp}(\lambda)$, $X \sim \text{Bin}(10, p)$, $X \sim N(m, \sigma^2)$)

3) necunoscută: $X \sim ?$

- în cazurile 2) și 3) parametrii necunoscuți sau distribuția necunoscută

↪ se estimează → teoria estimăției

↪ se testează → teste statistice

- Exemple de estimatori (funcții de selecție) sunt: media de selecție, dispersia de selecție, funcția de repartiție empirică (de selecție).

Fie x_1, \dots, x_n datele statistice pentru caracteristica cercetată X , notăm cu X_1, \dots, X_n variabilele de selecție corespunzătoare

► **media de selecție**

$$\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n)$$

► valoarea mediei de selecție

$$\bar{x}_n = \frac{1}{n} (x_1 + \dots + x_n)$$

► **varianța (dispersia) de selecție**

$$\tilde{S}_n^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$$

► valoarea varianței (dispersiei) de selecție

$$\tilde{s}_n^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x}_n)^2$$

► **abaterea standard de selecție**

$$\tilde{S}_n = \left(\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \right)^{\frac{1}{2}}$$

► valoarea abaterii standard de selecție

$$\tilde{s}_n = \left(\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x}_n)^2 \right)^{\frac{1}{2}}$$

► **funcția de repartiție empirică** $\hat{F}_n : \mathbb{R} \times \Omega \rightarrow \mathbb{R}$

$$\hat{F}_n(x) = \frac{\#\{i \in \{1, \dots, n\} : X_i \leq x\}}{n}, x \in \mathbb{R}$$

► valoarea funcției de repartiție empirice

$$\hat{F}_n(x) = \frac{\#\{i \in \{1, \dots, n\} : x_i \leq x\}}{n}, x \in \mathbb{R}$$

► $g(X_1, \dots, X_n)$ este **estimator nedeplasat** pentru parametrul necunoscut θ , dacă

$$E(g(X_1, \dots, X_n)) = \theta.$$

► $g(X_1, \dots, X_n)$ este **estimator consistent** pentru parametrul necunoscut θ , dacă

$$g(X_1, \dots, X_n) \xrightarrow{a.s.} \theta.$$

► Fie $g_1(X_1, \dots, X_n)$ și $g_2(X_1, \dots, X_n)$ estimatori nedeplasați pentru parametrul necunoscut θ . $g_1(X_1, \dots, X_n)$ este **mai eficient** decât $g_2(X_1, \dots, X_n)$, dacă $V(g_1) < V(g_2)$.

► Fie $\alpha \in (0, 1)$ nivelul de semnificație (probabilitatea de risc). **Cuantila de ordin α** pentru distribuția caracteristicii cercetate X este numărul $z_\alpha \in \mathbb{R}$ pentru care

$$P(X < z_\alpha) \leq \alpha \leq P(X \leq z_\alpha).$$

• dacă X este v.a. continuă, atunci: z_α este cuantilă de ordin $\alpha \iff P(X \leq z_\alpha) = \alpha \iff F_X(z_\alpha) = \alpha$

• $\alpha \cdot 100\%$ din valorile lui X sunt mai mici sau egale cu z_α

Exemple:

- 1) Media de selecție \bar{X}_n este estimator nedeplasat și consistent pentru **media teoretică** $E(X)$ a caracteristicii X .
- 2) Varianța (dispersia) de selecție \tilde{S}_n^2 este estimator nedeplasat și consistent pentru **varianța teoretică** $V(X)$ a caracteristicii X .
- 3) **Funcția de repartiție de selecție** calculată în $x \in \mathbb{R}$: $\hat{F}_n(x)$ este estimator nedeplasat și consistent pentru $F(x)$ valoarea funcției de repartiție teoretice în x .

Teste statistice

Fie x_1, \dots, x_n datele statistice pentru caracteristica cercetată X , notăm cu X_1, \dots, X_n variabilele de selecție corespunzătoare.

- Ipoteza statistică este o presupunere relativă la un parametru necunoscut θ
- Metoda de stabilire a veridicității unei ipoteze statistice se numește test (criteriu de verificare).
- Rezultatul testării se folosește apoi pentru luarea unor decizii (cum ar fi: eficiența unor medicamente, strategii de marketing, alegerea unui produs etc.).
- Se formulează ipoteza nulă H_0 și ipoteza alternativă H_1 , privind parametrul θ ; fie θ_0 o valoare dată

$$\text{I. } H_0 : \theta = \theta_0 \quad H_1 : \theta \neq \theta_0$$

$$\text{II. } H_0 : \theta \leq \theta_0 \quad H_1 : \theta > \theta_0$$

$$\text{III. } H_0 : \theta \geq \theta_0 \quad H_1 : \theta < \theta_0$$

Se dă $\alpha \in (0, 1)$ nivelul de semnificație (probabilitatea de risc). Formularea unui test revine la construirea unei regiuni critice $U \subset \mathbb{R}^n$ (pentru cazurile I, II, respectiv III) astfel încât

$$P((X_1, \dots, X_n) \in U | H_0) = \alpha$$

ceea ce este echivalent cu

$$P((X_1, \dots, X_n) \notin U | H_0) = 1 - \alpha$$

Concluzia testului:

$(x_1, \dots, x_n) \notin U \Rightarrow$ ipoteza H_0 este admisă

$(x_1, \dots, x_n) \in U \Rightarrow$ ipoteza H_0 este respinsă, în favoarea ipotezei H_1

- O colectivitate este testată în raport cu caracteristica X .
 - test pentru valoarea medie $E(X)$
 - ▷ când varianța teoretică $V(X)$ este cunoscută: testul lui Gauss (testul Z)
 - ▷ când varianța teoretică $V(X)$ este necunoscută: Student Test (testul T)
 - test pentru abaterea standard teoretică $\sqrt{V(X)}$ sau pentru varianța teoretică $V(X)$: testul χ^2
 - test asupra proporției (testul Gauss aproximativ)

Pași în efectuarea unui test statistic:

- Care parametru se testează? Care test este potrivit?
- Care este ipoteza nulă H_0 și care este ipoteza alternativă H_1 ?
- Care este nivelul de semnificație (probabilitatea de risc) α ?
- Calculul valorii estimatorului pe baza datelor statistice
- Concluzia testului

Test pentru media $m = E(X)$ caracteristicii cercetate X , când varianța $\sigma^2 = V(X)$ este cunoscută

- se dau $\alpha \in (0, 1)$, m_0 , σ , datele statistice x_1, \dots, x_n

- dacă $X \sim N(m, \sigma^2)$ sau $n > 30$ și X are o distribuție necunoscută, atunci $\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$

► folosind datele statistice x_1, \dots, x_n , se calculează $z = \frac{\bar{x}_n - m_0}{\frac{\sigma}{\sqrt{n}}}$

► cuantilele legii normale $N(0, 1)$:

$$z_{1-\frac{\alpha}{2}} = \text{norminv}(1 - \frac{\alpha}{2}, 0, 1), z_{1-\alpha} = \text{norminv}(1 - \alpha, 0, 1), z_\alpha = \text{norminv}(\alpha, 0, 1)$$

	I. $H_0: m = m_0$ $H_1: m \neq m_0$	II. $H_0: m \leq m_0$ $H_1: m > m_0$	III. $H_0: m \geq m_0$ $H_1: m < m_0$
Se acceptă H_0 dacă	$ z < z_{1-\frac{\alpha}{2}}$	$z < z_{1-\alpha}$	$z > z_\alpha$
Se respinge H_0 în favoarea lui H_1 , dacă	$ z \geq z_{1-\frac{\alpha}{2}}$	$z \geq z_{1-\alpha}$	$z \leq z_\alpha$

► în Octave/Matlab: *ztest*

► regiunea critică $U \subset \mathbb{R}^n$

$$\text{I. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \left| \frac{\bar{u}_n - m_0}{\frac{\sigma}{\sqrt{n}}} \right| \geq z_{1-\frac{\alpha}{2}} \right\}, \text{ unde } \bar{u}_n = \frac{1}{n} (u_1 + \dots + u_n)$$

$$\text{II. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{\bar{u}_n - m_0}{\frac{\sigma}{\sqrt{n}}} \geq z_{1-\alpha} \right\}$$

$$\text{III. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{\bar{u}_n - m_0}{\frac{\sigma}{\sqrt{n}}} \leq z_\alpha \right\}$$

Test pentru media $m = E(X)$ caracteristicii cercetate X , când varianța $\sigma^2 = V(X)$ este necunoscută

► se dau $\alpha \in (0, 1)$, m_0 , datele statistice x_1, \dots, x_n

► dacă $X \sim N(m, \sigma^2)$ sau $n > 30$ și X are o distribuție necunoscută, atunci $\frac{\bar{X}_n - m}{\frac{\bar{S}_n}{\sqrt{n}}} \sim \text{Student}(n-1)$

► folosind datele statistice x_1, \dots, x_n , se calculează $t = \frac{\bar{x}_n - m_0}{\frac{\bar{s}_n}{\sqrt{n}}}$

► cuantilele legii Student cu $n-1$ grade de libertate:

$$t_{1-\frac{\alpha}{2}} = \text{tinv}(1 - \frac{\alpha}{2}, n-1), t_{1-\alpha} = \text{tinv}(1 - \alpha, n-1), t_\alpha = \text{tinv}(\alpha, n-1)$$

	I. $H_0: m = m_0$ $H_1: m \neq m_0$	II. $H_0: m \leq m_0$ $H_1: m > m_0$	III. $H_0: m \geq m_0$ $H_1: m < m_0$
Se acceptă H_0 dacă	$ t < t_{1-\frac{\alpha}{2}}$	$t < t_{1-\alpha}$	$t > t_\alpha$
Se respinge H_0 în favoarea lui H_1 , dacă	$ t \geq t_{1-\frac{\alpha}{2}}$	$t \geq t_{1-\alpha}$	$t \leq t_\alpha$

► în Octave/Matlab: *ttest*

► regiunea critică $U \subset \mathbb{R}^n$

$$\text{I. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \left| \frac{\bar{u}_n - m_0}{\frac{\bar{\sigma}_n}{\sqrt{n}}} \right| \geq t_{1-\frac{\alpha}{2}} \right\}, \text{ unde } \bar{u}_n = \frac{1}{n} (u_1 + \dots + u_n), \bar{\sigma}_n = \left(\frac{1}{n-1} \sum_{k=1}^n (u_k - \bar{u}_n)^2 \right)^{\frac{1}{2}}$$

$$\text{II. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{\bar{u}_n - m_0}{\frac{\bar{\sigma}_n}{\sqrt{n}}} \geq t_{1-\alpha} \right\}$$

$$\text{III. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{\bar{u}_n - m_0}{\frac{\bar{\sigma}_n}{\sqrt{n}}} \leq t_\alpha \right\}$$

Test asupra proporției p pentru caracteristica $X \sim \text{Bernoulli}(p)$ (testul Gauss aproximativ)

► se dau $\alpha \in (0, 1)$, p_0 , datele statistice x_1, \dots, x_n

► dacă $X \sim \text{Bernoulli}(p)$ și $np(1-p) \geq 10$, atunci $\frac{\bar{X}_n - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1)$

► folosind datele statistice x_1, \dots, x_n , se calculează $z = \frac{\bar{x}_n - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$

► cuantilele legii normale $N(0, 1)$:

$$z_{1-\frac{\alpha}{2}} = \text{norminv}(1 - \frac{\alpha}{2}, 0, 1), z_{1-\alpha} = \text{norminv}(1 - \alpha, 0, 1), z_\alpha = \text{norminv}(\alpha, 0, 1)$$

	I. $H_0: p = p_0$ $H_1: p \neq p_0$	II. $H_0: p \leq p_0$ $H_1: p > p_0$	III. $H_0: p \geq p_0$ $H_1: p < p_0$
Se acceptă H_0 dacă	$ z < z_{1-\frac{\alpha}{2}}$	$z < z_{1-\alpha}$	$z > z_\alpha$
Se respinge H_0 în favoarea lui H_1 , dacă	$ z \geq z_{1-\frac{\alpha}{2}}$	$z \geq z_{1-\alpha}$	$z \leq z_\alpha$

► regiunea critică $U \subset \mathbb{R}^n$

$$\text{I. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \left| \frac{\bar{u}_n - p_0}{\frac{\sigma}{\sqrt{n}}} \right| \geq z_{1-\frac{\alpha}{2}} \right\}, \text{ unde } \bar{u}_n = \frac{1}{n} (u_1 + \dots + u_n)$$

$$\text{II. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{\bar{u}_n - p_0}{\frac{\sigma}{\sqrt{n}}} \geq z_{1-\alpha} \right\}$$

$$\text{III. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{\bar{u}_n - p_0}{\frac{\sigma}{\sqrt{n}}} \leq z_\alpha \right\}$$

Test pentru abaterea standard teoretică $\sigma = \sqrt{V(X)}$ a caracteristicii cercetate X / Test pentru varianța teoretică $\sigma^2 = V(X)$ a caracteristicii cercetate X

► se dau $\alpha \in (0, 1)$, σ_0 , datele statistice x_1, \dots, x_n

► dacă $X \sim N(m, \sigma^2)$, atunci $\frac{n-1}{\sigma^2} \tilde{S}_n^2 \sim \chi^2(n-1)$

► folosind datele statistice x_1, \dots, x_n , se calculează $c = \frac{n-1}{\sigma_0^2} \cdot \tilde{s}_n^2$

► cuantilele χ^2 (Chi-pătrat) cu $n-1$ grade de libertate:

$$c_{1-\frac{\alpha}{2}} = \text{chi2inv}(1 - \frac{\alpha}{2}, n-1), c_{\frac{\alpha}{2}} = \text{chi2inv}(\frac{\alpha}{2}, n-1), c_{1-\alpha} = \text{chi2inv}(1 - \alpha, n-1), c_\alpha = \text{chi2inv}(\alpha, n-1)$$

	I. $H_0: \sigma = \sigma_0$ $H_1: \sigma \neq \sigma_0$	II. $H_0: \sigma \leq \sigma_0$ $H_1: \sigma > \sigma_0$	III. $H_0: \sigma \geq \sigma_0$ $H_1: \sigma < \sigma_0$
Se acceptă H_0 , dacă	$c_{\frac{\alpha}{2}} < c < c_{1-\frac{\alpha}{2}}$	$c < c_{1-\alpha}$	$c > c_\alpha$
Se respinge H_0 în favoarea lui H_1 , dacă	$c \notin (c_{\frac{\alpha}{2}}, c_{1-\frac{\alpha}{2}})$	$c \geq c_{1-\alpha}$	$c \leq c_\alpha$

► în Octave/Matlab: *vartest*

► regiunea critică $U \subset \mathbb{R}^n$

$$\text{I. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{1}{\sigma_0^2} \sum_{k=1}^n (u_k - \bar{u}_n)^2 \notin (c_{\frac{\alpha}{2}}, c_{1-\frac{\alpha}{2}}) \right\}, \text{ unde } \bar{u}_n = \frac{1}{n} (u_1 + \dots + u_n)$$

$$\text{II. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{1}{\sigma_0^2} \sum_{k=1}^n (u_k - \bar{u}_n)^2 \geq c_{1-\alpha} \right\}$$

$$\text{III. } U = \left\{ (u_1, \dots, u_n) \in \mathbb{R}^n : \frac{1}{\sigma_0^2} \sum_{k=1}^n (u_k - \bar{u}_n)^2 \leq c_\alpha \right\}$$

Probleme:

1. Specificațiile unui anumit medicament indică faptul că fiecare comprimat conține în medie 2.4 g de substanță activă. 100 de comprimate alese la întâmplare din producție sunt analizate și se constată că ele conțin în medie 2.5 g de substanță activă cu o deviație standard de 0.2 g. Se poate spune că medicamentul respectă specificațiile (cu $\alpha = 0.01$)?

Soluție: $H_0: m = 2.4$ cu $H_1: m \neq 2.4$, testul Student.

2. Un manager este suspicios că un utilaj, care umple anumite cutii cu ceai, trebuie înlocuit cu unul mult mai precis. 121 de cutii cu ceai sunt cântărite. S-a obținut o medie de 196.6 g și o deviație standard de 2.09 g pentru acest eșantion.

a) Să se testeze dacă abaterea standard a utilajului este de 2 g.

b) Sunt datele suficiente pentru a concluziona, că utilajul trebuie reglat pentru că nu pune 200 g de ceai într-o cutie? ($\alpha = 0.01$)

Soluție: a) $H_0: \sigma = 2$ cu $H_1: \sigma \neq 2$, testul pentru abaterea standard

b) $H_0: m = 200$ cu $H_1: m \neq 200$, testul Student.

Intervale de încredere

Fie x_1, \dots, x_n datele statistice pentru caracteristica cercetată X , a cărei distribuție depinde de parametrul necunoscut θ ; notăm cu X_1, \dots, X_n variabilele de selecție corespunzătoare. Fie $\alpha \in (0, 1)$ nivelul de semnificație; $1 - \alpha$ se numește nivelul de încredere. Se caută doi estimatori $g_1(X_1, \dots, X_n)$ și $g_2(X_1, \dots, X_n)$ astfel încât

$$P(g_1(X_1, \dots, X_n) < \theta < g_2(X_1, \dots, X_n)) = 1 - \alpha \Leftrightarrow P(\theta \notin (g_1(X_1, \dots, X_n), g_2(X_1, \dots, X_n))) = \alpha$$

- ▶ $(g_1(X_1, \dots, X_n), g_2(X_1, \dots, X_n))$ se numește **interval de încredere pentru parametrul necunoscut θ**
- ▶ $(g_1(x_1, \dots, x_n), g_2(x_1, \dots, x_n))$ este valoarea intervalului de încredere pentru parametrul necunoscut θ
- ▶ $g_1(X_1, \dots, X_n)$ este limita inferioară a intervalului de încredere, valoarea sa este $g_1(x_1, \dots, x_n)$
- ▶ $g_2(X_1, \dots, X_n)$ este limita superioară a intervalului de încredere, valoarea sa este $g_2(x_1, \dots, x_n)$
- ▶ probabilitatea ca parametrul necunoscut θ să fie în intervalul $(g_1(X_1, \dots, X_n), g_2(X_1, \dots, X_n))$ este $1 - \alpha$ (nivelul de încredere)

Interval de încredere pentru media $m = E(X)$ caracteristicii cercetate X , când varianța $\sigma^2 = V(X)$ este cunoscută

▶ se dau $\alpha \in (0, 1)$, σ , datele statistice x_1, \dots, x_n

▶ dacă $X \sim N(m, \sigma^2)$ sau $n > 30$ și X are o distribuție necunoscută, atunci $\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$

▶ cuantilele legii normale $N(0, 1)$:

$$z_{1-\frac{\alpha}{2}} = \text{norminv}(1 - \frac{\alpha}{2}, 0, 1), z_{1-\alpha} = \text{norminv}(1 - \alpha, 0, 1), z_\alpha = \text{norminv}(\alpha, 0, 1)$$

• **interval de încredere bilateral:** să se indice doi estimatori $g_1(X_1, \dots, X_n)$ și $g_2(X_1, \dots, X_n)$ astfel încât pentru media teoretică $m = E(X)$ să avem:

$$P(g_1(X_1, \dots, X_n) < m < g_2(X_1, \dots, X_n)) = 1 - \alpha$$

$$\triangleright \text{din } P\left(\left|\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}}\right| < z_{1-\frac{\alpha}{2}}\right) = 1 - \alpha \text{ avem : } \left|\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}}\right| < z_{1-\frac{\alpha}{2}} \Leftrightarrow \bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_{1-\frac{\alpha}{2}} < m < \bar{X}_n + \frac{\sigma}{\sqrt{n}} \cdot z_{1-\frac{\alpha}{2}}$$

• **intervalul de încredere (bilateral)** pentru $m = E(X)$ (media teoretică) este $\left(\bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_{1-\frac{\alpha}{2}}, \bar{X}_n + \frac{\sigma}{\sqrt{n}} \cdot z_{1-\frac{\alpha}{2}}\right)$

$$\triangleright \text{se calculează valoarea intervalului de încredere } \left(\bar{x}_n - \frac{\sigma}{\sqrt{n}} \cdot z_{1-\frac{\alpha}{2}}, \bar{x}_n + \frac{\sigma}{\sqrt{n}} \cdot z_{1-\frac{\alpha}{2}}\right)$$

• **interval de încredere unilateral:** să se indice doi estimatori $g_1(X_1, \dots, X_n)$, $g_2(X_1, \dots, X_n)$ astfel încât pentru media teoretică

$m = E(X)$ să avem:

$$P(g_1(X_1, \dots, X_n) < m) = 1 - \alpha, \quad P(m < g_2(X_1, \dots, X_n)) = 1 - \alpha$$

$$\triangleright \text{din } P\left(\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} < z_{1-\alpha}\right) = 1 - \alpha \text{ avem : } \frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} < z_{1-\alpha} \Leftrightarrow m > \bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_{1-\alpha}$$

$$\bullet g_1(X_1, \dots, X_n) = \bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_{1-\alpha}; \text{ valoarea intervalului de încredere este } \left(\bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_{1-\alpha}, \infty\right)$$

$$\triangleright \text{din } P\left(\frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} > z_\alpha\right) = 1 - \alpha \text{ avem : } \frac{\bar{X}_n - m}{\frac{\sigma}{\sqrt{n}}} > z_\alpha \Leftrightarrow m < \bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_\alpha$$

$$\bullet g_2(X_1, \dots, X_n) = \bar{X}_n - \frac{\sigma}{\sqrt{n}} \cdot z_\alpha; \text{ valoarea intervalului de încredere este } \left(-\infty, \bar{x}_n - \frac{\sigma}{\sqrt{n}} \cdot z_\alpha\right)$$

Interval de încredere pentru media $m = E(X)$ caracteristicii cercetate X , când varianța $\sigma^2 = V(X)$ este necunoscută

► se dau $\alpha \in (0, 1)$, datele statistice x_1, \dots, x_n

► dacă $X \sim N(m, \sigma^2)$ sau $n > 30$ și X are o distribuție necunoscută, atunci $\frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}} \sim Student(n-1)$

► cuantilele legii $Student(n-1)$:

$$t_{1-\frac{\alpha}{2}} = norminv(1 - \frac{\alpha}{2}, 0, 1), t_{1-\alpha} = norminv(1 - \alpha, 0, 1), t_\alpha = norminv(\alpha, 0, 1)$$

• *interval de încredere bilateral*: să se indice doi estimatori $g_1(X_1, \dots, X_n)$ și $g_2(X_1, \dots, X_n)$ astfel încât pentru media teoretică $m = E(X)$ să avem:

$$P(g_1(X_1, \dots, X_n) < m < g_2(X_1, \dots, X_n)) = 1 - \alpha$$

$$\triangleright \text{din } P\left(\left|\frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}}\right| < t_{1-\frac{\alpha}{2}}\right) = 1 - \alpha \text{ avem : } \left|\frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}}\right| < t_{1-\frac{\alpha}{2}} \Leftrightarrow \bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2}} < m < \bar{X}_n + \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2}}$$

$$\bullet \text{intervalul de încredere (bilateral) pentru } m = E(X) \text{ (media teoretică) este } \left(\bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2}}, \bar{X}_n + \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2}}\right)$$

$$\triangleright \text{se calculează valoarea intervalului de încredere } \left(\bar{x}_n - \frac{\tilde{s}_n}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2}}, \bar{x}_n + \frac{\tilde{s}_n}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2}}\right)$$

• *interval de încredere unilateral*: să se indice doi estimatori $g_1(X_1, \dots, X_n)$, $g_2(X_1, \dots, X_n)$ astfel încât pentru media teoretică $m = E(X)$ să avem:

$$P(g_1(X_1, \dots, X_n) < m) = 1 - \alpha, \quad P(m < g_2(X_1, \dots, X_n)) = 1 - \alpha$$

$$\triangleright \text{din } P\left(\frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}} < t_{1-\alpha}\right) = 1 - \alpha \text{ avem : } \frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}} < t_{1-\alpha} \Leftrightarrow m > \bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\alpha}$$

$$\bullet g_1(X_1, \dots, X_n) = \bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\alpha}; \text{ valoarea intervalului de încredere este } \left(\bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_{1-\alpha}, \infty\right)$$

$$\triangleright \text{din } P\left(\frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}} > t_\alpha\right) = 1 - \alpha \text{ avem : } \frac{\bar{X}_n - m}{\frac{\tilde{S}_n}{\sqrt{n}}} > t_\alpha \Leftrightarrow m < \bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_\alpha$$

$$\bullet g_2(X_1, \dots, X_n) = \bar{X}_n - \frac{\tilde{S}_n}{\sqrt{n}} \cdot t_\alpha; \text{ valoarea intervalului de încredere este } \left(-\infty, \bar{x}_n - \frac{\tilde{s}_n}{\sqrt{n}} \cdot t_\alpha\right)$$

Interval de încredere pentru varianța (dispersia) $\sigma^2 = V(X)$ caracteristicii cercetate X

► se dau $\alpha \in (0, 1)$, datele statistice x_1, \dots, x_n

► dacă $X \sim N(m, \sigma^2)$, atunci $\frac{n-1}{\sigma^2} \tilde{S}_n^2 \sim \chi^2(n-1)$

► cuantilele χ^2 (Chi-pătrat) cu $n-1$ grade de libertate:

$$c_{1-\frac{\alpha}{2}} = \text{chi2inv}(1 - \frac{\alpha}{2}, n - 1), c_{\frac{\alpha}{2}} = \text{chi2inv}(\frac{\alpha}{2}, n - 1), c_{1-\alpha} = \text{chi2inv}(1 - \alpha, n - 1), c_{\alpha} = \text{chi2inv}(\alpha, n - 1)$$

• **interval de încredere bilateral:** să se indice doi estimatori $g_1(X_1, \dots, X_n)$ și $g_2(X_1, \dots, X_n)$ astfel încât pentru varianța teoretică $m = E(X)$ să avem:

$$P(g_1(X_1, \dots, X_n) < \sigma^2 < g_2(X_1, \dots, X_n)) = 1 - \alpha$$

$$\triangleright \text{din } P\left(c_{\frac{\alpha}{2}} < \frac{n-1}{\sigma^2} \cdot \tilde{S}_n^2 < c_{1-\frac{\alpha}{2}}\right) = 1 - \alpha \text{ avem : } c_{\frac{\alpha}{2}} < \frac{n-1}{\sigma^2} \cdot \tilde{S}_n^2 < c_{1-\frac{\alpha}{2}} \Leftrightarrow \frac{n-1}{c_{1-\frac{\alpha}{2}}} \cdot \tilde{S}_n^2 < \sigma^2 < \frac{n-1}{c_{\frac{\alpha}{2}}} \cdot \tilde{S}_n^2$$

• **intervalul de încredere (bilateral)** pentru $\sigma^2 = V(X)$ (varianța teoretică) este $\left(\frac{n-1}{c_{1-\frac{\alpha}{2}}} \cdot \tilde{S}_n^2 < \sigma^2 < \frac{n-1}{c_{\frac{\alpha}{2}}} \cdot \tilde{S}_n^2\right)$

$$\triangleright \text{se calculează valoarea intervalului de încredere } \left(\frac{n-1}{c_{1-\frac{\alpha}{2}}} \cdot \tilde{S}_n^2 < \sigma^2 < \frac{n-1}{c_{\frac{\alpha}{2}}} \cdot \tilde{S}_n^2\right)$$

• **interval de încredere unilateral:** să se indice doi estimatori $g_1(X_1, \dots, X_n), g_2(X_1, \dots, X_n)$ astfel încât pentru varianța teoretică $\sigma^2 = V(X)$ să avem:

$$P(g_1(X_1, \dots, X_n) < \sigma^2) = 1 - \alpha, \quad P(\sigma^2 < g_2(X_1, \dots, X_n)) = 1 - \alpha$$

$$\triangleright \text{din } P\left(\frac{n-1}{\sigma^2} \cdot \tilde{S}_n^2 < c_{1-\alpha}\right) = 1 - \alpha \text{ avem : } \frac{n-1}{\sigma^2} \cdot \tilde{S}_n^2 < c_{1-\alpha} \Leftrightarrow \sigma^2 > \frac{n-1}{c_{1-\alpha}} \cdot \tilde{S}_n^2$$

• $g_1(X_1, \dots, X_n) = \frac{n-1}{c_{1-\alpha}} \cdot \tilde{S}_n^2$; valoarea intervalului de încredere este $\left(\frac{n-1}{c_{1-\alpha}} \cdot \tilde{S}_n^2, \infty\right)$

$$\triangleright \text{din } P\left(\frac{n-1}{\sigma^2} \cdot \tilde{S}_n^2 > c_{\alpha}\right) = 1 - \alpha \text{ avem : } \frac{n-1}{\sigma^2} \cdot \tilde{S}_n^2 > c_{\alpha} \Leftrightarrow \sigma^2 < \frac{n-1}{c_{\alpha}} \cdot \tilde{S}_n^2$$

• $g_2(X_1, \dots, X_n) = \frac{n-1}{c_{\alpha}} \cdot \tilde{S}_n^2$; valoarea intervalului de încredere este $\left(-\infty, \frac{n-1}{c_{\alpha}} \cdot \tilde{S}_n^2\right)$

Metoda momentelor pentru estimarea parametrilor necunoscuți $\theta = (\theta_1, \dots, \theta_r)$ pentru distribuția caracteristicii cercetate X

de exemplu:

$X \sim \text{Exp}(\lambda)$ parametrul necunoscut: $\theta = \lambda$

$X \sim N(m, \sigma^2)$ parametri necunoscuți: $(\theta_1, \theta_2) = (m, \sigma)$

$X \sim \text{Unif}[a, b]$ parametri necunoscuți: $(\theta_1, \theta_2) = (a, b)$

Fie x_1, \dots, x_n datele statistice pentru caracteristica cercetată X și fie X_1, \dots, X_n variabilele de selecție corespunzătoare.

Se rezolvă sistemul

$$\begin{cases} E(X^k) = \frac{1}{n} \sum_{i=1}^n x_i^k, \\ k = \{1, \dots, r\} \end{cases}$$

cu necunoscutele $\theta_1, \dots, \theta_r$.

Soluția sistemului $\hat{\theta}_1, \dots, \hat{\theta}_r$ este estimatorul pentru parametrii necunoscuți ai distribuției caracteristicii X .

Exemplu: Folosind metoda momentelor, să se estimeze parametrul necunoscut $\theta := a$ pentru $X \sim \text{Unif}[0, a]$; se dau datele statistice: 0.1, 0.3, 0.9, 0.49, 0.12, 0.31, 0.98, 0.73, 0.13, 0.62.

Avem cazul: $r = 1$, calculăm $E(X) = \frac{a}{2}$, $n = 10$, $\bar{x}_n = 0.468$. Se rezolvă

$$E(X) = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \frac{a}{2} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Estimatorul pentru parametrul necunoscut a este

$$\hat{a}(X_1, \dots, X_n) = \frac{2}{n} \sum_{i=1}^n X_i,$$

unde X_1, \dots, X_n sunt variabilele de selecție. **Valoarea estimatorului** este

$$\hat{a}(x_1, \dots, x_n) = \frac{2}{n} \sum_{i=1}^n x_i = 0.936.$$

Parametrul necunoscut a este estimat cu valoarea 0.936.

► Este $\hat{a}(X_1, \dots, X_n)$ un estimator nedeplasat pentru parametrul a ?

Metoda verosimilității maxime pentru estimarea parametrului necunoscut θ al distribuției caracteristicii cercetate X

Fie x_1, \dots, x_n datele statistice pentru caracteristica cercetată X și fie X_1, \dots, X_n variabilele de selecție corespunzătoare. Notăm

$$L(x_1, \dots, x_n; \theta) = \begin{cases} P(X = x_1) \cdot \dots \cdot P(X = x_n), & \text{dacă } X \text{ e v.a. discretă} \\ f_X(x_1) \cdot \dots \cdot f_X(x_n), & \text{dacă } X \text{ e v.a. continuă.} \end{cases}$$

Aceasta este funcția de verosimilitate pentru parametrul θ și datele statistice x_1, \dots, x_n .

Metoda verosimilității maxime se bazează pe principiul că valoarea cea mai verosimilă (cea mai potrivită) a parametrului necunoscut θ este aceea pentru care funcția de verosimilitate $L(x_1, \dots, x_n; \theta)$ ia valoarea maximă:

$$(1) \quad L(x_1, \dots, x_n; \hat{\theta}) = \max_{\theta} L(x_1, \dots, x_n; \theta).$$

Se rezolvă sistemul $\frac{\partial L}{\partial \theta} = 0$ și se arată că $\frac{\partial^2 L}{\partial \theta^2} < 0$. Deseori este mai practic să se considere varianta transformată $\frac{\partial \ln L}{\partial \theta} = 0$ cu $\frac{\partial^2 \ln L}{\partial \theta^2} < 0$. În unele situații (1) se rezolvă prin alte metode.

Observație: Dacă distribuția caracteristicii cercetate depinde de k parametri necunoscuți $(\theta_1, \dots, \theta_k)$ atunci se rezolvă sistemul

$$\frac{\partial L}{\partial \theta_j} = 0, j = \overline{1, k} \text{ și se arată că matricea } \left(\frac{\partial^2 L}{\partial \theta_i \partial \theta_j} \right)_{1 \leq i \leq j \leq k} \text{ este negativ definită.}$$

Se poate lucra și cu varianta transformată:

$$\frac{\partial \ln L}{\partial \theta_j} = 0, j = \overline{1, k} \text{ și se arată că matricea } \left(\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right)_{1 \leq i \leq j \leq k} \text{ este negativ definită.}$$

O matrice M este negativ definită dacă $y^t M y < 0$ pentru orice $y \in \mathbb{R}^n \setminus \{0_n\}$.

Exemplu: Folosind metoda verosimilității maxime să se estimeze parametrul $\theta := p \in (0, 1)$ al distribuției Bernoulli,

$$X \sim \begin{pmatrix} 0 & 1 \\ 1-p & p \end{pmatrix}, \text{ cu datele statistice: } 0, 1, 1, 0, 0, 0, 1, 0, 1, 0.$$

$$\Rightarrow n = 10, x_1 = 0, x_2 = 1, x_3 = 1, x_4 = 0 \dots; P(X = x) = p^x (1-p)^{1-x}, x \in \{0, 1\}$$

$$\Rightarrow L(x_1, \dots, x_n; p) = P(X = x_1) \cdot \dots \cdot P(X = x_n) = p^{x_1 + \dots + x_n} (1-p)^{n - (x_1 + \dots + x_n)}$$

$$\Rightarrow \ln L(x_1, \dots, x_n; p) = (x_1 + \dots + x_n) \ln(p) + (n - (x_1 + \dots + x_n)) \ln(1-p)$$

$$\frac{\partial \ln L}{\partial p} = 0 \Rightarrow p = \frac{1}{n} (x_1 + \dots + x_n).$$

$$\text{Are loc: } \frac{\partial^2 \ln L}{\partial p^2} < 0.$$

Estimatorul de verosimilitate maximă pentru parametrul necunoscut p este

$$\hat{p}(X_1, \dots, X_n) = \frac{1}{n} (X_1 + \dots + X_n) = \bar{X}_n,$$

unde X_1, \dots, X_n sunt variabilele de selecție. **Valoarea estimată** este

$$\hat{p}(x_1, \dots, x_n) = \frac{1}{n} (x_1 + \dots + x_n) = \bar{x}_n = \frac{4}{10} = 0.4.$$

► Este $\hat{p}(X_1, \dots, X_n)$ un estimator nedeplasat pentru parametrul p ?