

Estimatori punctuali. Teorema limită centrală

Fie (X, f_θ) un **model statistic** asociat unei populații, adică X este o variabilă aleatoare ce cuantifică o caracteristică comună a indivizilor populației respective și a cărei distribuție depinde de un parametru necunoscut $\theta \in \Theta \subset \mathbb{R}^d$, $d \in \mathbb{N}^*$, iar f_θ este densitatea de probabilitate a variabilei aleatoare X în cazul în care X este o variabilă aleatoare continuă, respectiv $f_\theta(x) = P(X = x)$, $x \in \mathbb{R}$, dacă X este variabilă aleatoare discretă.

Vectorul θ poate să aibă o singură componentă, caz în care distribuția de probabilitate depinde de un singur parametru, de exemplu distribuția Bernoulli depinde de parametrul $p \in (0, 1)$, dar poate să fie și un vector cu mai multe componente, caz în care distribuția de probabilitate depinde de mai mulți parametri, de exemplu distribuția normală depinde de $m \in \mathbb{R}$ și $\sigma^2 > 0$, adică $\theta = (m, \sigma^2) \in \mathbb{R} \times (0, \infty)$.

Dacă X_1, X_2, \dots, X_n sunt variabile aleatoare independente identic distribuite (i.i.d.) ca și X , atunci vectorul (X_1, X_2, \dots, X_n) se numește **selecție aleatoare** de volum n asociată modelului statistic (X, f_θ) , iar dacă (x_1, x_2, \dots, x_n) este o realizare a selecției aleatoare (X_1, X_2, \dots, X_n) , caz în care se spune că x_1, x_2, \dots, x_n sunt valori de selecție ale variabilei aleatoare X , atunci orice funcție

$$\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n),$$

de variabile x_1, x_2, \dots, x_n , se numește **estimator punctual** al parametrului θ . Cu alte cuvinte, un estimator punctual al parametrului θ este o realizare a variabilei aleatoare $\hat{\theta}(X_1, X_2, \dots, X_n)$.

Tipuri de estimatori:

Un estimator $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ al parametrului θ se numește:

- **estimator consistent** dacă

$$\lim_{n \rightarrow \infty} P\left(|\hat{\theta}(X_1, X_2, \dots, X_n) - \theta| > \varepsilon\right) = 0, \text{ pentru orice } \varepsilon > 0,$$

adică distribuția variabilei $\hat{\theta}(X_1, X_2, \dots, X_n)$ este din ce în ce mai concentrată în jurul parametrului θ pe măsură ce volumul selecției aleatoare crește;

- **estimator nedeplasat (centrat)** dacă

$$M\left(\hat{\theta}(X_1, X_2, \dots, X_n)\right) = \theta.$$

Fie $\hat{\theta}_1, \hat{\theta}_2$ doi estimatori nedeplasați ai parametrului θ . Spunem că $\hat{\theta}_1$ este **mai eficient** decât $\hat{\theta}_2$ dacă are loc inegalitatea:

$$\sigma^2 \left(\hat{\theta}_1(X_1, X_2, \dots, X_n) \right) < \sigma^2 \left(\hat{\theta}_2(X_1, X_2, \dots, X_n) \right),$$

ceea ce este echivalent cu

$$\frac{\sigma^2 \left(\hat{\theta}_1(X_1, X_2, \dots, X_n) \right)}{\sigma^2 \left(\hat{\theta}_2(X_1, X_2, \dots, X_n) \right)} < 1.$$

Cu alte cuvinte, cu cât varianța (dispersia) variabilei $\hat{\theta}(X_1, X_2, \dots, X_n)$ este mai mică, cu atât estimatorul $\hat{\theta}$ este mai eficient.

- Determinarea mediei și dispersiei pentru variabila medie aritmetică \overline{X} .

Dacă (X_1, X_2, \dots, X_n) este o selecție aleatoare de volum n asociată unui model statistic (X, f_θ) , cu $M(X) = m$ și $\sigma^2(X) = \sigma^2$, atunci variabila medie aritmetică

$$\overline{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

are media m și dispersia $D^2 = \frac{\sigma^2}{n}$.

În plus, dacă variabilele aleatoare i.i.d. X_1, X_2, \dots, X_n au o **distribuție normală**,

$$X_i \sim N(m, \sigma^2), i = \overline{1, n},$$

atunci media lor aritmetică \overline{X} are de asemenea o **distribuție normală** de medie m și dispersie $D^2 = \frac{\sigma^2}{n}$, adică

$$\overline{X} \sim N(m, D^2 = \sigma^2/n). \quad (1)$$

- **Teorema limită centrală** (enunț, consecințe)

Dacă (X_n) este un șir de variabile aleatoare independente și identic distribuite având media comună m și abaterea standard σ , iar

$$\overline{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n}$$

este șirul variabilelor medie aritmetică, atunci pentru n suficient de mare ($n \geq 30$) distribuția de probabilitate a variabilelor \overline{X}_n este **aproximativ normală** de medie m și dispersie $D^2 = \sigma^2/n$. Notăm

$$\overline{X}_n \sim ApN(m, D^2 = \sigma^2/n).$$

În plus, pentru n suficient de mare, avem

$$S_n = X_1 + X_2 + \cdots + X_n \sim ApN(nm, D^2 = n\sigma^2). \quad (2)$$

- **Estimarea mediei:** Media de selecție, definită prin

$$\overline{x} = \frac{x_1 + x_2 + \cdots + x_n}{n},$$

este un estimator nedeplasat al mediei $m = M(X)$, adică

$$\widehat{m} = \overline{x}$$

este un estimator nedeplasat al parametrului m .

- **Estimarea varianței (dispersiei):** Dispersia de selecție, definită prin

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

este un estimator nedeplasat al dispersiei $\sigma^2 = \sigma^2(X)$.

Estimatorul verosimilității maxime (maximum likelihood estimator):

Fie (X, f_θ) un model statistic și x_1, x_2, \dots, x_n valori de selecție ale lui X (se mai spune că x_1, x_2, \dots, x_n este un eșantion de volum n). **Funcția de verosimilitate** (likelihood function) asociată eșantionului x_1, x_2, \dots, x_n este $L : \Theta \subset \mathbb{R}^d \rightarrow \mathbb{R}$,

$$L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i) = f_\theta(x_1) \cdot f_\theta(x_2) \cdots f_\theta(x_n).$$

Punctul de maxim global al funcției de verosimilitate (dacă există) se numește **estimatorul verosimilității maxime** al parametrului

$$\theta = (\theta_1, \theta_2, \dots, \theta_d) \in \Theta.$$

Cu alte cuvinte, estimatorul verosimilității maxime al lui θ (dacă acesta există) este

$$\hat{\theta} = \operatorname{argmax} L(\theta),$$

unde prin $\operatorname{argmax} L(\theta)$ se înțelege argumentul care maximizează funcția L .

Probleme rezolvate

Exemplul 1. Cererea de memorie pentru o aplicație, ca proporție din memoria ce poate fi alocată de un utilizator, este o variabilă aleatoare X ce are densitatea de probabilitate

$$f(x) = \begin{cases} (\theta + 1)x^\theta, & 0 < x < 1, \\ 0, & \text{în rest.} \end{cases}$$

- a) Să se determine media teoretică $M(X)$ a variabilei aleatoare X și apoi un estimator nedeplasat al parametrului θ în funcție de media de selecție \bar{x} a unei selecții aleatoare de volum n .
- b) Să se determine un estimator nedeplasat al parametrului θ din selecția următoare:

$$0.2, 0.4, 0.5, 0.7, 0.8, 0.9, 0.9, 0.6, 0.6, 0.4,$$

rezultată în urma rulării aplicației cu diferite date de intrare.

Rezolvare: a) Mai întâi calculăm media teoretică:

$$M(X) = \int_{-\infty}^{\infty} xf(x)dx = \int_0^1 (\theta + 1)x^{\theta+1}dx = (\theta + 1) \frac{x^{\theta+2}}{\theta + 2} \Big|_0^1 = \frac{\theta + 1}{\theta + 2}.$$

Dacă $m = M(X)$ și \bar{x} este media de selecție a unui eșantion x_1, x_2, \dots, x_n , atunci din egalitatea impusă

$$\hat{m} = \bar{x}$$

se determină un estimator nedeplasat al parametrului θ :

$$\frac{\hat{\theta} + 1}{\hat{\theta} + 2} = \bar{x} \quad \Leftrightarrow \quad \hat{\theta} = \frac{2\bar{x} - 1}{1 - \bar{x}}.$$

- b) Pentru valorile înregistrate deducem că $\bar{x} = 0.6$, deci un estimator nedeplasat al parametrului θ este $\hat{\theta} = \frac{2\bar{x} - 1}{1 - \bar{x}} = 0.5$.

□

Exemplul 2. Fie \mathcal{P} populația formată dintr-un tip de circuite. Caracteristica ce dorim să o investigăm este durata de viață a acestor circuite, știind că aceasta este exponențial distribuită, cu parametrul $\theta > 0$ necunoscut.

Măsurând timpul de viață (în ani) a 10 circuite, se obțin valorile:

0.8830, 1.9651, 1.9189, 4.8448, 0.9208 3.4377, 1.7162, 4.2327, 5.9435, 8.3128.

Să se determine estimatorul verosimilității maxime al parametrului θ (adică pentru media duratei de viață a acestui tip de circuite).

Soluție. Densitatea de probabilitate a distribuției exponențiale este

$$f_{\theta}(x) = \begin{cases} 0, & \text{dacă } x < 0, \\ \frac{1}{\theta}e^{-\frac{x}{\theta}}, & \text{dacă } x \geq 0. \end{cases}$$

Astfel, funcția de verosimilitate este

$$L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = \prod_{i=1}^n \left(\frac{1}{\theta} e^{-\frac{x_i}{\theta}} \right) = \frac{1}{\theta^n} e^{-\frac{\sum_{i=1}^n x_i}{\theta}}.$$

Notăm cu

$$\ell(\theta) = \ln(L(\theta)).$$

Funcțiile ℓ și L au aceleași puncte de extrem și de aceeași natură, deoarece funcția $x \mapsto \ln x$ este strict crescătoare. Pentru simplificarea calculelor, vom determina punctul de maxim global (dacă acesta există) pentru ℓ , iar acesta va fi punct de maxim global și pentru L . Avem

$$\ell(\theta) = \ln L(\theta) = -n \ln \theta - \frac{\sum_{i=1}^n x_i}{\theta},$$

deci

$$\ell'(\theta) = -\frac{n}{\theta} + \frac{\sum_{i=1}^n x_i}{\theta^2}.$$

Rezolvând ecuația

$$\ell'(\theta) = 0,$$

în raport cu θ , obținem punctul

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

care este punct de maxim global pentru ℓ , deci și pentru L . În concluzie, estimatorul verosimilității maxime al parametrului θ al distribuției exponențiale este media de selecție, $\hat{\theta} = \bar{x}$. În cazul exemplului dat, estimatorul verosimilității maxime al mediei de viață a circuitelor este:

$$\hat{\theta} = \bar{x} = \frac{x_1 + x_2 + \dots + x_{10}}{10} = 3.4175.$$

Observație: Valorile de selecție alese mai sus au fost generate simulând o variabilă aleatoare $X \sim \text{Exp}(\theta = 3.5)$,

> rexp(10, 1/3.5) # în limbajul de programare R, rate=1/θ,

deci estimatorul verosimilității maxime $\hat{\theta} = 3.4175$ este destul de bun.

□

Exemplul 3. Un simulator al distribuției Bernoulli

$$X : \begin{pmatrix} 1 & 0 \\ p & 1-p \end{pmatrix},$$

de parametru $p \in (0, 1)$ necunoscut, generează stringul de biți:

1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0.

Să se determine estimatorul verosimilității maxime al parametrului p pe baza eșantionului de biți generați.

Soluție: În cazul de față parametrul necunoscut este $\theta = p$. Cum X este o variabilă aleatoare discretă, avem

$$f_p(x) = P(X = x) = \begin{cases} p, & \text{dacă } x = 1, \\ 1 - p, & \text{dacă } x = 0. \end{cases}$$

Funcția de verosimilitate asociată eșantionului de biți generați este

$$L(p) = L(p; x_1, x_2, \dots, x_{25}) = f_p(x_1) \cdot f_p(x_2) \cdots f_p(x_{25}) = p^{12}(1-p)^{13},$$

căci stringul are 12 biți de 1 și 13 biți de 0. Se logaritmează și se determină punctul de maxim global al funcției

$$\ell(p) = \ln L(p) = 12 \ln(p) + 13 \ln(1-p).$$

Cum

$$\ell'(p) = \frac{12}{p} - \frac{13}{1-p} = \frac{12(1-p) - 13p}{p(1-p)} = \frac{12 - 25p}{p(1-p)},$$

ecuația

$$\ell'(p) = 0$$

are soluția $p = \frac{12}{25}$. Se arată că $p = \frac{12}{25}$ este un punct de maxim al funcției ℓ , deci estimatorul verosimilității maxime pentru parametrul p al distribuției Bernoulli, dedus din stringul de biți generați mai sus, este

$$\hat{p} = \frac{12}{25} = 0.48,$$

adică este egal cu numărul biților de 1 din string supra numărul total de biți. Acest estimator al lui p este de fapt probabilitatea intuitivă de a obține bitul 1: numărul cazurilor favorabile supra numărul cazurilor posibile din string.

□

Probleme propuse

1. Numărul de accesări ale `http://cs.upt.ro` într-o oră este o variabilă aleatoare Poisson X de parametru $\lambda > 0$ necunoscut. În 5 ore consecutive numărul de accesări înregistrate a fost

10, 12, 18, 15, 15.

Să se estimeze parametrul λ din aceste date. Să se determine estimatorul verosimilității maxime al lui λ .

2. Există ipoteza că timpul de așteptare T , pentru a primi acces la un server internet, are densitatea de probabilitate

$$f(x) = \begin{cases} \theta x^{\theta-1}, & 0 < x < 1, \\ 0, & \text{în rest,} \end{cases}$$

unde θ este un parametru pozitiv necunoscut. Pentru a estima parametrul θ s-au făcut 5 observații care au condus la următorii timpi de așteptare (în secunde):

0.4, 0.7, 0.2, 0.6, 0.1.

Să se calculeze media teoretică a timpului de așteptare și apoi un estimator nedeplasat al parametrului θ . Să se determine estimatorul verosimilității maxime pentru parametrul θ pe baza datelor din eșantion.

3. Deduceți estimatorul verosimilității maxime al parametrului p al distribuției $\text{Bin}(n, p)$, unde n este un număr natural nenul **cunoscut**.

4. Șirul de valori de mai jos

3.9, 2.4, 2.1, 3.6, 3.0

a fost generat apelând de 5 ori consecutiv generatorul `2 + (b - 2) * urand()`; Să se estimeze parametrul b .

Indicație: Numerele date sunt valori de observație asupra unei variabile aleatoare $X \sim \text{Unif}[2, b]$. Media unei variabile aleatoare $X \sim \text{Unif}[a, b]$ este $M(X) = (a + b)/2$. Parametrul b se poate calcula determinând un estimator nedeplasat al mediei $m = M(X)$, $\hat{m} = \bar{x}$.

5. Să se determine estimatorul verosimilității maxime pentru parametrul θ al distribuției uniforme

$$f_{\theta}(x) = \begin{cases} \frac{1}{\theta}, & x \in [0, \theta], \\ 0, & \text{în rest,} \end{cases}$$

pe baza valorilor de observație x_1, x_2, \dots, x_n .

Indicație: Funcția de verosimilitate este

$$L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = \frac{1}{\theta^n}, \text{ dacă } \theta \geq \max\{x_1, x_2, \dots, x_n\}.$$

Funcția $L : [\max\{x_i : i = \overline{1, n}\}, \infty) \rightarrow \mathbb{R}$ are derivata $L'(\theta) = -n/\theta^{n+1} < 0$, deci este strict descrescătoare și maximul global îl atinge în $\max\{x_i : i = \overline{1, n}\}$. Deci estimatorul verosimilității maxime este $\hat{\theta} = \max\{x_i : i = \overline{1, n}\}$.

6. Fie X_i , $i = \overline{1, 100}$, variabile aleatoare i.i.d. de distribuție Bernoulli($p = 0.2$). Să se precizeze distribuția de probabilitate a variabilei medie aritmetică

$$\overline{X} = \frac{1}{100} \sum_{i=1}^{100} X_i$$

și să se estimeze probabilitatea $P(\overline{X} > 0.3)$.

7. Un DVD conține o zonă neînregistrată de 350 MB. Să se estimeze probabilitatea ca această zonă de memorie să fie suficientă pentru a copia 300 de fotografii ce au dimensiuni independente, de medie 1 MB și abatere standard 0.5 MB.