

Curs 11

Statistică inferențială

Statistica inferențială se concentrază asupra luării deciziilor despre caracteristicile unei populații bazate pe informațiile obținute prin studiul datelor uneia sau mai multor selecții aleatoare și pe calcule probabilistice.

Statistica inferențială se împarte în două mari domenii:

- estimarea parametrilor și
- ipoteze statistice.

Metodele de inferență se bazează pe respectarea caracterului aleator al selecției.

Fie X o v. a. cu *legea de probabilitate* $f(x, \theta)$, $x \in \mathbb{R}$, $\theta \in \Theta \subset \mathbb{R}$, care poate fi funcția de frecvență, în cazul discret, respectiv densitatea de probabilitate, în cazul continuu.

Definiția 11.0.1 *Mulțimea legilor de probabilitate $f(x, \theta)$ care conțin parametrul necunoscut θ se numește model probabilistic.*

În practică există numeroase astfel de familii de legi de probabilitate care pot fi alese ca modele probabilistice: normală, exponentială, Poisson etc. Alegerea unei astfel de familii, ca model al unui fenomen particular real se face pe baza experienței anterioare în studiul fenomenelor asemănătoare sau în urma unui studiu preliminar al rezultatelor experienței sau observației.

În momentul în care am ales o lege de probabilitate de o anumită formă, incertitudinea legată de rezultatul particular al experimentului s-a transferat în incertitudinea legată de valorile parametrului (parametrilor).

Definiția 11.0.2 Repartiția selecției $X = (X_1, X_2, \dots, X_n)$ este definită ca repartiția comună a variabilelor de selecție X_1, X_2, \dots, X_n . Dacă x_1, x_2, \dots, x_n sunt valorile luate de variabilele de selecție X_1, X_2, \dots, X_n atunci legea de probabilitate a selecției este notată prin

$$f(x_1, x_2, \dots, x_n; \theta), (x_1, x_2, \dots, x_n) \in \mathbb{R}^n.$$

Repartiția selecției depinde de θ și încorporează atât selecția cât și modelul probabilistic.

Cea mai folosită formă de selecție este selecția aleatoare și este bazată pe ideea experimentului aleator.

Definiția 11.0.3 Vom spune că (X_1, X_2, \dots, X_n) este o selecție aleatoare asupra v. a. X care are legea de probabilitate $f(x, \theta)$ dacă X_1, X_2, \dots, X_n sunt independente și identic distribuite (i.i.d.) ca și X .

În cazul selecției aleatoare legea de probabilitate a selecției este

$$f(x_1, x_2, \dots, x_n; \theta) = \prod_{j=1}^n f(x_j, \theta).$$

O selecție aleatoare poate fi constituită prin repetarea unui experiment aleator de n ori. Realizarea selecției aleatoare (datele obținute în urma selecției) se notează cu (x_1, x_2, \dots, x_n) și mulțimea tuturor realizărilor definește spațiul observațiilor.

Definiția 11.0.4 Modelul probabilist $f(x, \theta)$ împreună cu selecția aleatoare

$$X = (X_1, X_2, \dots, X_n)$$

definesc modelul statistic.

Modelul statistic împreună cu datele observate impun considerarea următoarelor întrebări:

1. datele observate sunt consistente cu modelul statistic postulat?
2. presupunând că modelul statistic postulat este consistent cu datele observate, ce putem spune despre parametrii necunoscuți $\theta \in \Theta$?
 - a) putem descrește incertitudinea asupra lui θ prin reducerea spațiului parametrilor Θ la Θ_0 unde Θ_0 este o submulțime a lui Θ ? (estimație prin intervale de încredere)
 - b) putem descrește incertitudinea asupra lui θ prin alegerea unei valori particulare $\hat{\theta}$ din Θ ca având cea mai reprezentativă valoare a lui θ ? (estimație punctuală)
 - c) putem considera întrebarea: θ aparține unei submulțimi Θ_0 a lui Θ ? (verificarea ipotezelor statistice)
3. presupunând că a fost aleasă o valoare particulară $\hat{\theta}$ reprezentativă a lui θ , ce putem spune despre observațiile viitoare? (predicție)

Definiția 11.0.5 Dată o populație de volum N , se numește statistică a acelei populații orice mărime semnificativă calculată cu ajutorul datelor obținute dintr-un eșantion aleator al acelei populații.

Exemple de statistici frecvent folosite: media de selecție statistică, dispersia de selecție statistică, dispersia de selecție modificată.

Definiția 11.0.6 Pentru un parametru al repartiției, se numește estimator al aceluia parametru orice statistică care aproximează acel parametru.

Dacă considerăm selecția aleatoare $X = (X_1, X_2, \dots, X_n)$ cu datele obținute în urma selecției (x_1, x_2, \dots, x_n) , vom nota estimatorul parametrului θ cu $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ și atunci putem scrie

$$\forall \varepsilon > 0 : \lim_{n \rightarrow \infty} P \left(\left| \hat{\theta}(x_1, x_2, \dots, x_n) - \theta \right| > \varepsilon \right) = 0$$

Cu alte cuvinte, pentru $\varepsilon > 0$ dat, pentru valori mari ale lui n este foarte puțin probabil ca $\hat{\theta}(x_1, x_2, \dots, x_n)$ să ia valori în afara intervalului $[\theta - \varepsilon, \theta + \varepsilon]$, adică este foarte puțin probabil ca numărul $\hat{\theta}(x_1, x_2, \dots, x_n)$ să fie în afara intervalului $[\theta - \varepsilon, \theta + \varepsilon]$. În aceste condiții, după un număr de n experiențe, considerăm pe $\hat{\theta}(x_1, x_2, \dots, x_n)$ ca o aproximație bună pentru θ . Este posibil să ne înșelăm, dar probabilitatea de a ne înșela este mică, pentru n mare. Statistica nu ne oferă răspunsuri sigure ci doar aproximații în care putem avea un grad mai mic sau mai mare de încredere. Se acceptă acele aproximații în care avem un grad mai mare de încredere.

Se utilizează estimări punctuale sau estimări cu ajutorul intervalelor de încredere. Estimatorul punctual atribuie o valoare parametrului. Intervalele de încredere generează un interval în care parametrul aparține cu un anumit nivel de încredere.

Exemplul 11.0.7 Considerăm populația diametrelor unor bile de rulmenți dintr-un lot. Ca statistică a acestei populații se poate lua media de selecție \bar{x} a unui șir (x_1, x_2, \dots, x_n) format de diametrele unui eșantion de bile din acel lot. Considerând alt lot, media de selecție \bar{x} va avea altă valoare.

Exemplul 11.0.8 O societate de telefoane poate considera populația tuturor abonaților. Un eșantion se poate obține luând pe sărite din lista abonaților, de exemplu, din 25 în 25. Presupunem că societatea este interesată în estimarea unui parametru caracteristic; de exemplu, gradul de satisfacție g al abonaților relativ la serviciile oferite de societate. Se cere o apreciere cu note de la 1 la 5. Este solicitat răspunsul abonaților din eșantionul ales și se obține un estimator al lui g .

Dacă, de exemplu, θ reprezintă media atunci putem considera ca statistică a mediei, $\theta = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$. Fie x_1, x_2, \dots, x_n sunt valorile luate de X_1, X_2, \dots, X_n atunci estimatorul lui θ este $\hat{\theta} = \bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$. În condițiile în care măsurătorile sunt realizate cu acuratețe și precizie rezonabile, atunci X_1, X_2, \dots, X_n se pot considera independente și identic repartizate.

Vom prezenta exemple de modele statistice care apar adesea în practică.

11.1 Modele statistice

Modelul Bernoulli

Presupunem că X_1, X_2, \dots, X_n este o selecție aleatoare a unei populații care urmează o distribuție Bernoulli cu parametrul $\theta \in [0, 1]$ necunoscut. Știm că distribuția Bernoulli poate fi utilizată la observarea unui articol, rezultat al unui proces industrial, caz în care $X_i = 1$ indică faptul că articolul i este bun și $X_i = 0$ dacă articolul este defect. În studii medicale rezultatul $X_i = 1$ indică faptul că tratamentul aplicat pacientului i a avut succes iar $X_i = 0$ în caz contrar. În aceste cazuri vrem să știm valoarea lui θ .

Spațiul parametrului θ este $[0, 1]$. Dacă x_1, x_2, \dots, x_n sunt valorile luate de X_1, X_2, \dots, X_n atunci *funcția de frecvență* pentru X_i este

$$f(x_i, \theta) = \theta^{x_i}(1 - \theta)^{1-x_i},$$

și reprezintă $P(X = x_i)$, θ este parametrul necunoscut, iar legea de probabilitate a selecției este dată de

$$\prod_{i=1}^n f(x_i, \theta) = \prod_{i=1}^n \theta^{x_i}(1 - \theta)^{1-x_i} = \theta^{n\bar{x}}(1 - \theta)^{n(1-\bar{x})}.$$

Modelul Poisson

Presupunem că X_1, X_2, \dots, X_n este o selecție aleatoare a unei populații care urmează o distribuție Poisson cu parametrul $\theta \in \mathbb{R}_+$ necunoscut. Știm că distribuția Repartiția Poisson a apărut odată cu studiul variabilelor aleatoare cu un număr foarte mare de valori. Exemple clasice de astfel de variabile aleatoare sunt cele legate de **evenimentele care apar rar** într-un interval de timp. În aceste cazuri vrem să știm valoarea lui θ .

Spațiul parametrului θ este \mathbb{R}_+ . Dacă x_1, x_2, \dots, x_n sunt valorile luate de X_1, X_2, \dots, X_n atunci *funcția de frecvență* pentru X_i este

$$f(x_i, \theta) = \frac{\theta^{x_i}}{x_i!} e^{-\theta},$$

și reprezintă $P(X = x_i)$, θ este parametrul necunoscut, iar legea de probabilitate a selecției este dată de

$$\prod_{i=1}^n f(x_i, \theta) = \prod_{i=1}^n \frac{\theta^{x_i}}{x_i!} e^{-\theta} = \frac{\theta^{n\bar{x}}}{x_1! x_2! \dots x_n!} e^{-n\theta}.$$

Modelul normal

Presupunem că X_1, X_2, \dots, X_n este o selecție aleatoare a unei populații care urmează o distribuție normală $N(\mu, \sigma)$ cu $\theta = (\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+$ necunoscuți, unde $\mathbb{R}^+ = (0, \infty)$. De exemplu, putem avea observații asupra înălțimii în centimetri a unei populații și intuim că este rezonabil să presupunem că distribuția înălțimii este normală cu media și dispersia necunoscute. Dacă x_1, x_2, \dots, x_n sunt valorile luate de X_1, X_2, \dots, X_n atunci densitatea de probabilitate a unui eșantion este dată de

$$\prod_{i=1}^n f(x_i; \mu, \sigma^2) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

Deoarece

$$\sum_{i=1}^n (x_i - \mu)^2 = n(\bar{x} - \mu)^2 + \sum_{i=1}^n (x_i - \bar{x})^2,$$

unde $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ este media de selecție și $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ este dispersia de selecție modificată (deviația standard), obținem

$$\prod_{i=1}^n f(x_i; \mu, \sigma^2) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{n}{2\sigma^2} (\bar{x} - \mu)^2 - \frac{n-1}{2\sigma^2} s^2\right). \quad (11.1)$$

Ne putem pune întrebarea dacă presupunerea că înălțimea populației urmează o distribuție normală este corectă. Această ipoteză trebuie verificată. Procedurile care duc la verificarea unor astfel de ipoteze se numesc teste statistice.

Modelul exponențial

Presupunem că timpul de viață a unui aparat este distribuit $Exp[\theta]$ unde $\theta \in (0, \infty)$ este necunoscut. Presupunem că X_1, X_2, \dots, X_n este o selecție aleatoare a unei populații care urmează o distribuție exponențială. Dacă x_1, x_2, \dots, x_n sunt valorile luate de X_1, X_2, \dots, X_n atunci densitatea de probabilitate a unui eșantion este dată de

$$\prod_{i=1}^n f(x_i; \theta) = \theta^n \exp\{-n\bar{x}\theta\}.$$

Modelul multinomial

În capitolul II am introdus distribuția multinomială. Aceasta apare în aplicații în care în urma unui experiment se pot realiza k evenimente diferite independente A_1, A_2, \dots, A_k . Aceste evenimente se produc cu probabilitățile p_1, p_2, \dots, p_k , $\sum_{j=1}^k p_j = 1$, $0 \leq p_j \leq 1$.

Facem presupunerea că, de exemplu, $k = 3$. În acest caz nu cunoaștem parametrii (p_1, p_2, p_3) . Dacă x_1, x_2, x_3 sunt valorile luate de X_1, X_2, X_3 atunci densitatea de probabilitate a unui eșantion este dată de

$$\prod_{i=1}^3 f(x_i; p_1, p_2, p_3) = p_1^{x_1} p_2^{x_2} p_3^{x_3}.$$

Construcția histogramei și încercarea de a intui modelul statistic sunt metode folosite de statistician în încercarea de a studia distribuția populației. Deseori nu este clar ce model statistic trebuie folosit. Mai mult, dacă am intuit modelul statistic, nu știm nimic despre valoarea parametrilor care intra în definiția modelului statistic. De exemplu, știm că înălțimea urmează o distribuție normală, dar nu știm media și dispersia. Folosind intuiția am putea presupune că media poate fi aproximată cu \bar{x} . Pentru a justifica această alegere avem nevoie să studiem estimatorii punctuali.

11.2 Estimatori punctuali

Estimarea punctuală este valoarea atribuită unui parametru pe baza statisticii construite din eșantion.

Definiția 11.2.1 *Se consideră o populație de volum N și un parametru θ al acestei populații. Fie X_1, X_2, \dots, X_n este o selecție aleatoare dintr-un eșantion reprezentativ ($n \ll N$) al populației care ia valorile x_1, x_2, \dots, x_n ; un estimator punctual $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$ a lui θ se numește estimator nedepășat sau absolut corect dacă $M[\hat{\theta}] = \theta$.*

În practică se recomandă $n \geq 30$ și $n \leq \frac{N}{3}$.

Definiția 11.2.2 *Dacă $\lim_{n \rightarrow \infty} M[\hat{\theta}] = \theta$ și $\lim_{n \rightarrow \infty} D[\hat{\theta}] = 0$ atunci $\hat{\theta}$ este un estimator corect sau depășat al lui θ (în cazul estimatorului consistent limita trebuie să fie oricât de mică atunci când n crește).*

Teorema 11.2.3 *Dacă statistica $\hat{\theta}(x_1, x_2, \dots, x_n)$ este un estimator corect a lui θ atunci este un estimator al lui θ , adică pentru orice $\varepsilon > 0$ avem*

$$\lim_{n \rightarrow \infty} P \left(\left| \hat{\theta}(x_1, x_2, \dots, x_n) - \theta \right| > \varepsilon \right) = 0.$$

Demonstrație. Conform inegalității lui Cebâșev,

$$P \left(\left| \hat{\theta}(x_1, x_2, \dots, x_n) - M[\hat{\theta}] \right| > \varepsilon \right) \leq \frac{D[\hat{\theta}]}{\varepsilon^2},$$

dar $\lim_{n \rightarrow \infty} M[\hat{\theta}] = \theta$ și $\lim_{n \rightarrow \infty} D[\hat{\theta}] = 0$ de unde rezultă concluzia. \square

Arătăm că statisticile introduse anterior sunt estimatori ale parametrilor corespunzători.

Teorema 11.2.4 *Fie x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale variabilelor X_1, X_2, \dots, X_n (v. a. independente și identic distribuite ca și X). Notăm*

$m = M[X]$ media teoretică,

$\sigma^2 = D[X]$ dispersia teoretică,

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ media de selecție,

$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ dispersia de selecție și

$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ dispersia de selecție modificată.

Atunci media de selecție \bar{X} și dispersia de selecție S^2 sunt estimatori punctuali pentru media, respectiv dispersia teoretică. În plus,

1. media de selecție \bar{X} este un estimator absolut corect al lui m ;
2. dispersia de selecție $\bar{\sigma}^2$ este un estimator corect al lui σ^2 .
3. dispersia de selecție modificată S^2 este un estimator absolut corect al lui σ^2 .

Demonstrație. 1. Trebuie verificate condițiile din definiția estimatorului absolut corect.

$$M[\bar{X}] = \frac{1}{n} M\left[\sum_{i=1}^n X_i\right] = \frac{1}{n} n M[X] = M[X].$$

$$D[\bar{X}] = D\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} D\left[\sum_{i=1}^n X_i\right] = \frac{1}{n^2} \sum_{i=1}^n D[X_i] = \frac{\sigma^2}{n}.$$

2.

$$\begin{aligned} M[\bar{\sigma}^2] &= \frac{1}{n} M\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] = \\ &= \frac{1}{n} M\left[\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2\right] = \\ &= \frac{1}{n} M\left[\sum_{i=1}^n X_i^2 - n\bar{X}^2\right] = M\left[\frac{1}{n} \sum_{i=1}^n X_i^2\right] - M[\bar{X}^2]. \end{aligned}$$

Dar

$$M[\bar{X}^2] = \frac{1}{n^2} M\left[\left(\sum_{i=1}^n X_i\right)^2\right] = \frac{1}{n^2} M\left[\sum_{i=1}^n X_i^2\right] + \frac{2}{n^2} \sum_{i < j} M[X_i X_j].$$

Apoi $M[X_i^2] = M[X^2]$ și din ipoteza de independență

$$M[X_i X_j] = M[X_i] M[X_j] = (M[X])^2,$$

rezultă

$$M[\bar{X}^2] = \frac{1}{n^2} M[X^2] + \frac{n-1}{n} (M[X])^2.$$

Deci

$$\begin{aligned} M[\bar{\sigma}^2] &= \frac{1}{n} M[X^2] - \left(\frac{1}{n^2} M[X^2] + \frac{n-1}{n} (M[X])^2\right) = \\ &= \frac{n-1}{n} (M[X^2] - (M[X])^2) = \frac{n-1}{n} \sigma^2 \rightarrow \sigma^2. \end{aligned}$$

Rezultă că $\bar{\sigma}^2$ nu estimează absolut corect dispersia v. a. X .

3. Reținem că $M[\bar{\sigma}^2] = \frac{n-1}{n} \sigma^2$.

Rezultă că $M[S^2] = M\left[\frac{n}{n-1} \bar{\sigma}^2\right] = \sigma^2$. Se poate demonstra că

$$D[S^2] = \frac{1}{n} \left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^4 - \frac{n-3}{n-1} D[X^2] \right), \quad (11.2)$$

de unde rezultă că $D[S^2] \rightarrow 0$. □

Observația 11.2.5 Deși dispersia de selecție S^2 este un estimator nedeplasat a lui σ^2 , S este un estimator deplasat a lui σ . Pentru valori mari $n \geq 10$ se poate considera că $\hat{\sigma} = \left(1 + \frac{1}{4(n-1)}\right) S$ este un estimator nedeplasat a lui σ .

Observația 11.2.6 În practică se folosește S^2 în locul lui $\overline{\sigma^2}$ deoarece dă rezultate mai bune după cum ne arată Teorema 11.2.4. Totuși formula (11.2) ne spune că pentru n suficient de mare și statistica $\overline{\sigma^2}$ poate fi folosită ca estimator al dispersiei v.a. X .

Următorul tabel sintetizează relația dintre parametri necunoscuți și estimatorii punctuali asociați. Astfel μ notează media, σ^2 este dispersia, p este probabilitatea repartiției binomiale, $p = \frac{k}{n}$, unde k este numărul realizărilor într-un eșantion de volum n , $\mu_1 - \mu_2$ estimează diferența dintre mediile a două eșantioane independente a două statistici aleatoare, $p_1 - p_2$ estimează diferența dintre proporțiile a două eșantioane independente a două statistici aleatoare.

Param. nec.	Statistica	Estim. punctual
μ	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
σ^2	$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
p	$\hat{p} = \frac{1}{n} X$	$\hat{p} = \frac{1}{n} x$
$\mu_1 - \mu_2$	$\bar{X}_1 - \bar{X}_2 =$ $= \frac{1}{n_1} \sum_{i=1}^{n_1} X_{1i} - \frac{1}{n_2} \sum_{i=1}^{n_2} X_{2i}$	$\bar{x}_1 - \bar{x}_2 =$ $= \frac{1}{n_1} \sum_{i=1}^{n_1} x_{1i} - \frac{1}{n_2} \sum_{i=1}^{n_2} x_{2i}$
$p_1 - p_2$	$\hat{p}_1 - \hat{p}_2 = \frac{1}{n_1} X_1 - \frac{1}{n_2} X_2$	$\hat{p}_1 - \hat{p}_2 = \frac{1}{n_1} x_1 - \frac{1}{n_2} x_2$

Observația 11.2.7 Dacă $X \in \text{Poiss}[\lambda]$, atunci $M[X] = D[X] = \lambda$ și atunci \bar{x} și s^2 sunt estimatorii nedeplasați pentru λ . Se pune întrebarea pe care îl vom alege. Cum dispersia este o măsură a împrăștierei, intuiția sugerează să-l alegem pe acela care are cea mai mică dispersie și aceasta deoarece el are o repartiție mai concentrată în jurul lui λ .

Dintre doi estimatori nedeplasați $\hat{\theta}$ și $\hat{\theta}_1$ pentru același parametru θ , se consideră că este mai precis cel care are dispersie mai mică.

Teorema 11.2.8 Fie X o v. a. cu media m și abaterea medie pătratică σ . Dacă x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale lui X (v. a. independente la fel distribuite ca și X), atunci pentru n suficient de mare se poate considera că $\bar{X} \in N\left(m, \frac{\sigma}{\sqrt{n}}\right)$.

Demonstrație. Conform teoremei limită centrală avem că

$$S_n = X_1 + X_2 + \dots + X_n \Rightarrow S_n \in N(nm, \sigma\sqrt{n}).$$

Deoarece

$$M\left[\frac{1}{n}S_n\right] = \frac{1}{n}nm = m, \quad D\left[\frac{1}{n}S_n\right] = \frac{1}{n^2}n\sigma^2 = \frac{\sigma^2}{n},$$

rezultă că $\frac{1}{n}S_n \in N\left(m, \frac{\sigma}{\sqrt{n}}\right)$. □

Corolarul 11.2.9 În condițiile Teoremei 11.2.8 pentru orice $a < b$ avem

$$P(a \leq \bar{X} < b) \approx \Phi\left(\frac{b-m}{\frac{\sigma}{\sqrt{n}}}\right) - \Phi\left(\frac{a-m}{\frac{\sigma}{\sqrt{n}}}\right). \quad (11.3)$$

Statisticienii recomandă folosirea formulei (11.3) pentru populații statistice de volum N și eşantioane de volum n unde $n \geq 30$ și $n \leq \frac{N}{3}$. Dacă $n < 30$ formula este utilă dacă populația inițială nu este departe de a fi normal distribuită. Pentru eşantioane mici statisticienii propun o corecție care să înlocuiască $\frac{\sigma}{\sqrt{n}}$ prin $\frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$.

Exercițiul 1 Considerăm greutatea populației de 350 studenți dintr-o facultate s-a constatat că media este $m = 70$ și abaterea media pătratică este $\sigma = 10$.

a) Să se determine probabilitatea ca un student luat la întâmplare să cântărească între 65 și 70 kg?

b) Să se determine probabilitatea ca media maselor să fie între 65 și 75 kg.

c) Să se determine probabilitatea ca extrăgând un eşantion de 36 studenți din acea populație, media eşantionului să fie cuprinsă între 65 și 75 kg.

d) Extrăgând un eşantion de 100 de studenți care este probabilitatea ca media maselor să fie sub 65 kg?

Rezolvare. a) Probabilitatea cerută nu poate fi calculată deoarece $X =$ greutatea populației nu este repartizată normal.

La fel și la b).

c) $n = 36$, $m = 70$, $\sigma_{\bar{X}} = \frac{10}{\sqrt{36}} = 1.6667$,

$$\begin{aligned} P(65 \leq \bar{X} \leq 75) &= \Phi\left(\frac{75 - 70}{1.6667}\right) - \Phi\left(\frac{65 - 70}{1.6667}\right) = \\ &= 0.9986 - 0.0014 = 0.9972. \end{aligned}$$

d) $n = 100$, $m = 70$, $\sigma_{\bar{X}} = \frac{10}{\sqrt{100}} = 1$,

$$P(\bar{X} \leq 65) = \Phi\left(\frac{65 - 70}{1}\right) = 0.0000003. \quad \diamond$$

11.3 Metoda momentelor

Idea metodei este de a egala momentele distribuției cu momentele statistice și de aici obținem estimarea parametrilor. Momentele populației vor depinde de parametrii necunoscuți.

De exemplu momentul inițial de ordin întâi este media, iar media statistică este momentul statistic de ordin întâi. Egalăm acestea și obținem estimarea punctuală $\hat{\mu} = \bar{X}$. Deci media statistică este un estimator punctual al mediei populației.

Fie x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale variabilelor X_1, X_2, \dots, X_n care au o distribuție exponențială cu parametru λ . Avem un singur parametru de estimat. Știm că $M[X] = \frac{1}{\lambda}$. De aici rezultă că $\frac{1}{\lambda} = \bar{X} \Rightarrow \hat{\lambda} = \frac{1}{\bar{x}}$ este un estimator punctual al parametrului λ obținut cu metoda momentelor.

Exemplul 11.3.1 Ca exemplu, presupunem că este testat timpul de viață a unui modul electronic folosit în industria automobilelor. Timpul de viață urmează o distribuție exponențială. S-au făcut teste și s-au obținut următoarele date: $x_1 = 11.96$, $x_2 = 5.03$, $x_3 = 67.40$, $x_4 = 16.07$, $x_5 = 31.50$, $x_6 = 7.73$, $x_7 = 11.10$, $x_8 = 22.38$. Deoarece $\bar{x} = 21.65$, estimarea punctuală a parametrului λ obținut cu metoda momentelor este $\hat{\lambda} = \frac{1}{21.65} = 0.046189$.

Fie x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale variabilelor X_1, X_2, \dots, X_n care au o distribuție normală de parametri μ și σ . În acest caz avem $M[X] = \mu$, $M[X^2] = \mu^2 + \sigma^2$.

Rezultă că

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \mu^2 + \sigma^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i,$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^2 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Reamintim că acest estimator pentru σ^2 este un estimator corect sau deplasat.

11.4 Metoda verosimilității maxime

Una din cele mai bune metode de a obține un estimator punctual pentru un parametrul θ este metoda verosimilității maxime.

Metoda verosimilității maxime a fost dezvoltată în 1920 de matematicianul englez R. A Fisher și extinsă de Cramer, Rao și Wald. Este cea mai folosită metodă de estimare și joacă un rol important în verificarea ipotezelor statistice.

Fie v. a. X cu densitatea de probabilitate $f(x, \theta)$, unde θ este parametrul necunoscut. Considerăm o selecție X_1, X_2, \dots, X_n i.i.d cu v. a. X și fie x_1, x_2, \dots, x_n valorile variabilelor de selecție. Atunci *funcția de verosimilitate maximă* corespunzătoare selecției este

$$L(x_1, x_2, \dots, x_n; \theta) = \prod_{j=1}^n f(x_j, \theta).$$

Această funcție se utilizează numai dacă avem un singur parametru necunoscut θ . Estimatorul de verosimilitate maximă este acea valoare a lui θ care maximizează funcția $L(x_1, x_2, \dots, x_n; \theta)$. În cazul unei v. a. discrete interpretarea funcției este

$$L(x_1, x_2, \dots, x_n; \theta) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n),$$

adică reprezintă tocmai probabilitatea de a obține în urma selecției valorile x_1, x_2, \dots, x_n .

Metoda verosimilității maxime constă în următorul principiu (axiomă): valoarea "cea mai verosimilă" (cea mai potrivită în acest sens!) a parametrului este aceea pentru care funcția $L(x_1, x_2, \dots, x_n; \theta)$ este maximă. Aceasta implică θ punct în care

$$\frac{\partial L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} = 0. \quad (11.4)$$

Ecuția (11.4) în practică se dovedește a fi dificilă. De aceea se folosește observația: $L(x_1, x_2, \dots, x_n; \theta)$ este maximă dacă și numai dacă $\ln L(x_1, x_2, \dots, x_n; \theta)$ este maximă.

Introducem funcțiile: *logaritmul funcției de verosimilitate* $\ln L(x_1, x_2, \dots, x_n; \theta)$ și *funcția scor*

$$\frac{\partial L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} = s(x_1, x_2, \dots, x_n; \theta).$$

Acestea încorporează aceeași informație ca și $L(x_1, x_2, \dots, x_n; \theta)$. Dacă avem oricare dintre ele putem să o obținem pe cealaltă.

Definiția 11.4.1 $t_n(x_1, x_2, \dots, x_n) : S \rightarrow \Theta$ se numește estimatie de verosimilitate maximă pentru parametrul θ dacă este punct de maxim pentru funcția de verosimilitate, adică

$$\ln L(x_1, x_2, \dots, x_n; t_n) \geq \ln L(x_1, x_2, \dots, x_n; \theta)$$

pentru toți $\theta \in \Theta$.

În cele ce urmează presupunem că $f(x, \theta)$ este derivabilă până la ordinul doi inclusiv în raport cu θ . Atunci există $\frac{\partial L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta}$ și estimatia de verosimilitate maximă este soluție a ecuației

$$\frac{\partial L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} = 0 \quad (11.5)$$

sau a ecuației

$$\frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} = 0. \quad (11.6)$$

Fie această soluție θ' . Pentru a fi estimator de verosimilitate maximă trebuie să satisfacă și condiția

$$\left. \frac{\partial^2 \ln L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta^2} \right|_{\theta=\theta'} < 0.$$

Definiția 11.4.2 Ecuațiile (11.5) și (11.6) se numesc ecuații de verosimilitate.

Exemplul 11.4.3 Fie x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale variabilelor X_1, X_2, \dots, X_n care au o distribuție binomială, i. i. r. cu v. a. $X \in \text{Bernoulli}[\theta]$. Atunci

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

θ este parametrul care trebuie estimat. Funcția de verosimilitate maximă este, pentru selecția făcută,

$$\begin{aligned} L(x_1, x_2, \dots, x_n; \theta) &= \theta^{x_1} (1 - \theta)^{1-x_1} \theta^{x_2} (1 - \theta)^{1-x_2} \dots \theta^{x_n} (1 - \theta)^{1-x_n} = \\ &= \prod_{j=1}^n \theta^{x_j} (1 - \theta)^{1-x_j} = \theta^{n\bar{x}} (1 - \theta)^{n-n\bar{x}}. \end{aligned}$$

Observăm că θ minimizează $L(x_1, x_2, \dots, x_n; \theta)$ dacă minimizează și $\ln L(x_1, x_2, \dots, x_n; \theta)$. Atunci

$$\begin{aligned} \ln L(x_1, x_2, \dots, x_n; \theta) &= n\bar{x} \ln \theta + (n - n\bar{x}) \ln (1 - \theta). \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} &= \frac{n\bar{x}}{\theta} - \frac{n - n\bar{x}}{1 - \theta} \\ \frac{n\bar{x}}{\theta} - \frac{n - n\bar{x}}{1 - \theta} &= 0 \Rightarrow \hat{\theta} = \bar{x}. \end{aligned}$$

De aici rezultă că estimatorul de verosimilitate maximă pentru p este

$$\hat{\theta} = \frac{1}{n} \sum_{j=1}^n x_j. \quad \diamond$$

Exemplul 11.4.4 Fie x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale variabilelor de selecție X_1, X_2, \dots, X_n care au o distribuție normală, i. i. d. cu v. a. $X \in N(m, \sigma)$.

a) Dacă σ^2 este cunoscută și media m este necunoscută cu $\Theta = \mathbb{R}$, atunci funcția de verosimilitate maximă este

$$L(x_1, x_2, \dots, x_n; m) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - m)^2 \right\}.$$

Deoarece funcția $\ln L$ este strict crescătoare în L și valoarea lui m care maximizează funcția $\ln L(x_1, x_2, \dots, x_n; m)$ maximizează și $L(x_1, x_2, \dots, x_n; m)$, este mai ușor de lucrat cu logaritmul lui L decât cu L . Avem

$$\ln L(x_1, x_2, \dots, x_n; m) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - m)^2.$$

$$\text{Dar } \frac{\partial \ln L(x_1, x_2, \dots, x_n; m)}{\partial m} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - m),$$

$$\text{deci } \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - m) = 0 \Rightarrow \hat{m} = \bar{x}.$$

Observăm că $\frac{\partial^2 \ln L(x_1, x_2, \dots, x_n; m)}{\partial m^2} = -\frac{n}{\sigma^2} \leq 0, \forall m \in \mathbb{R}$, deci $\hat{m} = \bar{x}$ este punct de maxim. Astfel $t_n(x_1, x_2, \dots, x_n) = \bar{x}$ este punct de verosimilitate maximă.

b) Dacă m este cunoscută și $\sigma^2 \in \Theta = (0, \infty)$ este necunoscut folosim forma funcției de verosimilitate maximă dedusă în (11.1) și obținem

$$\begin{aligned} L(x_1, x_2, \dots, x_n; m, \sigma^2) &= \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (\bar{x} - m)^2 - \frac{n-1}{2\sigma^2} s^2 \right\} = \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (\bar{x} - m)^2 \right\} \exp \left\{ -\frac{n-1}{2\sigma^2} s^2 \right\}. \end{aligned}$$

Fixăm $m = \bar{x}$ și obținem

$$\begin{aligned} L(x_1, x_2, \dots, x_n; \sigma^2) &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{n-1}{2\sigma^2} s^2 \right\} \Rightarrow \\ \ln L(x_1, x_2, \dots, x_n; \sigma^2) &= -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{n-1}{2\sigma^2} s^2. \end{aligned}$$

$$\begin{aligned} \frac{\partial \ln L(x_1, x_2, \dots, x_n; \sigma^2)}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{n-1}{2\sigma^4} s^2, \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \sigma^2)}{\partial \sigma^2} &= 0 \Rightarrow \hat{\sigma}^2 = \frac{n-1}{n} s^2. \end{aligned}$$

Dar

$$\begin{aligned} \frac{d^2 \ln L}{d(\sigma^2)^2}(x_1, x_2, \dots, x_n; \hat{\sigma}^2) &= \frac{n}{2\sigma^4} - \frac{n-1}{\sigma^6} s^2 \Big|_{\sigma^2 = \frac{n-1}{n} s^2} = \\ &= -\frac{1}{2} \frac{n^3}{s^4 (n-1)^2} \leq 0. \end{aligned}$$

Rezultă că $\hat{\sigma}^2$ este estimatorul de verosimilitate maximă.

Dacă avem de estimat mai mulți parametri $1, 2, \dots, p$, cu X cu densitatea de probabilitate $f(x, \theta_1, \theta_2, \dots, \theta_p)$, atunci în mod analog cu cazul unui singur parametru, aplicăm principiul verosimilității maxime. Considerăm o selecție X_1, X_2, \dots, X_n i.i.d cu v. a. X și fie x_1, x_2, \dots, x_n valorile variabilelor de selecție. Atunci *funcția de verosimilitate maximă* corespunzătoare selecției este

$$L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_p) = \prod_{j=1}^n f(x_j, \theta_1, \theta_2, \dots, \theta_p).$$

Se aleg pentru parametri acele valori care maximizează funcția de verosimilitate

$$L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_p)$$

sau ceea ce este același lucru acele valori care maximizează $\ln L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_p)$.

Aceasta implică:

$$\begin{cases} \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_p)}{\partial \theta_1} = 0 \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_p)}{\partial \theta_2} = 0 \\ \dots \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_p)}{\partial \theta_p} = 0 \end{cases} \quad (11.7)$$

și diferențiala de ordin doi, care este o formă pătratică, calculată în punctul $(\theta_1, \theta_2, \dots, \theta_p)$, soluție a sistemului (11.7), negativ definită.

Exemplul 11.4.5 Fie x_1, x_2, \dots, x_n ($n \geq 2$) o selecție de valori ale variabilelor de selecție X_1, X_2, \dots, X_n care au o distribuție normală, i. i. d. cu v. a. $X \in N(m, \sigma)$. Dacă m și σ^2 sunt necunoscute, $\Theta = \{(m, \sigma^2) | m \in \mathbb{R}, \sigma^2 \in (0, \infty)\}$, atunci estimatorul de verosimilitate maximă $(\hat{m}, \hat{\sigma}^2)$ se obține rezolvând sistemul

$$\begin{cases} \frac{\partial \ln L(x_1, x_2, \dots, x_n; m, \sigma^2)}{\partial m} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - m) = 0 \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; m, \sigma^2)}{\partial \sigma^2} = 0 \end{cases}$$

adică

$$\begin{cases} \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - m) = 0 \\ -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (\bar{x} - m)^2 + \frac{n-1}{2(\sigma^2)^2} s^2 = 0 \end{cases}.$$

Obținem

$$\begin{cases} \hat{m} = \bar{x} \\ \hat{\sigma}^2 = \frac{n-1}{n} s^2 \end{cases}$$

Calculăm

$$\begin{aligned} d^2 \ln L(x_1, x_2, \dots, x_n; m, \sigma^2) &= \frac{\partial^2 \ln L(x_1, \dots, x_n; m, \sigma^2)}{\partial m^2} dm^2 + \\ + 2 \frac{\partial^2 \ln L(x_1, \dots, x_n; m, \sigma^2)}{\partial m \partial \sigma^2} dm d\sigma^2 + \frac{\partial^2 \ln L(x_1, \dots, x_n; m, \sigma^2)}{(\partial \sigma^2)^2} (d\sigma^2)^2, \\ d^2 \ln L(x_1, x_2, \dots, x_n; \hat{m}, \hat{\sigma}^2) &= -\frac{1}{\sigma^2} dm^2 - \frac{n^3}{2(n-1)^2 (s^2)^2} (d\sigma^2)^2 \end{aligned}$$

care este o formă pătratică, evident, negativ definită. De aici rezultă că $\left(\bar{x}, \frac{n-1}{n}s^2\right)$ este unicul estimator de verosimilitate maximă pentru (m, σ^2) . \diamond

Exemplul 11.4.6 Presupunem că timpul de viață a unui aparat este distribuit $Exp[\theta]$ unde $\theta \in (0, \infty)$ este necunoscut. Bazându-ne pe o selecție (x_1, x_2, \dots, x_n) de valori ale variabilelor X_1, X_2, \dots, X_n care au o distribuție exponențială obținem funcția de verosimilitate

$$\begin{aligned} L(x_1, x_2, \dots, x_n; \theta) &= \theta^n \exp\{-n\bar{x}\theta\} \Rightarrow \\ \ln L(x_1, x_2, \dots, x_n; \theta) &= n \ln \theta - n\bar{x}\theta \Rightarrow \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} &= \frac{n}{\theta} - n\bar{x} \Rightarrow \frac{n}{\theta} - n\bar{x} = 0 \Rightarrow \hat{\theta} = \frac{1}{\bar{x}}. \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 \ln L}{\partial \theta^2}(x_1, x_2, \dots, x_n; \theta) &= -\frac{n}{\theta^2} \Rightarrow \\ \frac{\partial^2 \ln L}{\partial \theta^2}(x_1, x_2, \dots, x_n; \frac{1}{\bar{x}}) &= -n(\bar{x})^2 \leq 0 \Rightarrow \end{aligned}$$

$\hat{\theta} = \frac{1}{\bar{x}}$ este estimatorul de verosimilitate maximă a lui θ . \diamond