

Curs 13: Introducere în Statistică. Estimatori

1.1 Problematika statisticii

Dezvoltarea tehnologiei a condus la generarea unui volum imens de date. Pe WEB se generează date în format text, imagine sau alt format multimedia. A crescut cantitatea de date numerice generate în experimentele din fizica energiilor înalte, în astronomie, explorarea spațiului, biologie etc. Pe de altă parte, se generează date în banking, telekom și în tranzacțiile mediului de afaceri. Aceste date de volum uriaș ascund informație care trebuie extrasă și utilizată pentru a facilita avansul în domeniile respective.

Există mai multe domenii care dezvoltă tehnici și proceduri de înregistrare a datelor de analiză și extragere a informației și a cunoștințelor din date. Pe de o parte, *statistica* are o istorie îndelungată în această direcție, iar pe de altă parte *machine learning*, *data mining* și, mai nou, *data science* sunt domenii noi care au apărut și s-au dezvoltat pe măsură ce a avansat tehnologia sistemelor de calcul.

Machine learning este știința care se ocupă cu proiectarea, analiza, implementarea și aplicațiile programelor ce învață din experiență sau își îmbunătățesc performanțele automat, prin experiență. Acest domeniu are o intersecție mare cu statistica. În timp ce *statistica* are ca scop formularea inferențelor pe baza informației extrase dintr-un eșantion de date, *machine learning* încorporează aspecte adiționale relativ la analiza algoritmilor ce pot fi folosiți pentru a capta, stoca, indexa, extrage și combina aceste date în scopul de a îndeplini sarcini greu de realizat prin mijloace algoritmice clasice.

Printre domeniile care folosesc instrumentele oferite de *machine learning* amintim: *computer vision* (recunoașterea fețelor, detectarea și localizarea obiectelor), *information retrieval* (extragerea informației din documentele indexate de către motoarele de căutare, identificarea topicilor din *feed*-uri, regăsirea imaginilor), navigare autonomă, controlul roboților, traducere automată (Google translate), bioinformatică și multe altele. *Machine learning* este unul din domeniile cu cea mai ridicată rată de dezvoltare. Companii mari, ca Google, Yahoo, IBM, Microsoft, investesc enorm în programe de cercetare-dezvoltare în *machine learning*.

Studiile experimentale de laborator, efectuate în diverse domenii din inginerie, fizică, chimie sau experimentele pe calculator din diverse domenii ale științei și tehnologiei, constând din simulări ale unor sisteme simple sau complexe, implică investigații statistice

ale datelor observate, măsurate sau simulate. Investigarea statistică constă în a studia o caracteristică comună a unei mulțimi de elemente de aceeași natură, numită *populație*. Elementele unei populații se numesc, generic, *indivizi*. Scopul investigației statistice este de a extrage informații despre caracteristica populației, investigând doar un eșantion constând din n indivizi, selectați la întâmplare. Numărul n al indivizilor din eșantion se numește *volumul eșantionului*.

Caracteristica comună a indivizilor populației este cuantificată de o variabilă aleatoare X , pentru care fie nu se cunoaște distribuția de probabilitate (densitatea de probabilitate f sau funcția de repartiție F), fie se cunoaște doar parțial, în sensul că se cunoaște tipul de distribuție de probabilitate a caracteristicii investigate, dar densitatea de probabilitate f_θ a variabilei aleatoare continue X sau distribuția de probabilitate p_θ , când X este variabilă aleatoare discretă ($p_\theta(x) = P(X = x)$), depinde de un parametru necunoscut $\theta \in \Theta \subseteq \mathbb{R}^d$, $d \geq 1$.

Observând sau măsurând caracteristica indivizilor dintr-un eșantion, se obține un șir de valori x_1, x_2, \dots, x_n , interpretate ca valori de observație asupra variabilei aleatoare X . Din acestea "se estimează" parametrii de interes, cum ar fi media caracteristicii investigate, dispersia sau parametrii necunoscuți, de care depinde distribuția de probabilitate.

Mai precis, statistica dezvoltă metode bazate pe rezultate din teoria probabilităților, care permit estimarea parametrului necunoscut $\theta \in \mathbb{R}^d$ al modelului probabilist f_θ sau p_θ .

În procesul de observare sau măsurare a caracteristicii indivizilor dintr-un eșantion se consideră că rezultatul investigării unui individ este independent de cel al investigării celorlalți. De aceea valorile înregistrate, x_1, x_2, \dots, x_n , sunt interpretate ca valori de observație asupra unui șir de variabile aleatoare X_1, X_2, \dots, X_n , independente și identic distribuite ca variabila aleatoare X , ce modelează caracteristica investigată. Practic, înțelegem prin X_k ca fiind caracteristica individului k din eșantion, $k = \overline{1, n}$. Cele n variabile aleatoare, având aceeași distribuție ca și X , au atât media $m = M(X_k)$, cât și dispersia $\sigma^2 = \sigma^2(X_k)$ egale cu cele ale lui X .

Considerăm \mathcal{P} o populație supusă investigării statistice din punctul de vedere al unei caracteristici X , ce ia valori discrete sau continue. Perechea (X, f_θ) , unde X este o variabilă aleatoare reală de densitate f_θ , dacă X este continuă, respectiv $f_\theta(x) = p_\theta(x)$, adică $f_\theta(x) = P(X = x)$, dacă X este discretă, se numește *model statistic*.

De exemplu, populația poate fi un anumit tip de chip-uri și caracteristica pe care dorim s-o investigăm este durata de viață. În general, durata de viață a dispozitivelor și circuitelor este exponențial distribuită. Astfel, modelul statistic al caracteristicii durată de viață este (X, f_θ) cu $X \sim \text{Exp}(\theta)$, θ fiind un parametru necunoscut. Înregistrând durata de viață a n chipuri selectate la întâmplare din producția dintr-o anumită perioadă, se va estima parametrul θ , care se știe că reprezintă media variabilei aleatoare X a modelului exponențial. Având un estimator al lui θ , firma producătoare poate stabili garanția pe care o dă pentru buna funcționare a tipului respectiv de chip-uri.

Nu întotdeauna populația constă din obiecte fizice, palpabile. De exemplu, pentru a deduce distribuția de probabilitate a intervalului dintre două pachete de informație pe un canal de comunicație, populația investigată constă din astfel de intervale. În acest caz eșantionul nu se alege apriori și apoi să se facă măsurătorile, ci se observă direct într-o

perioadă dată, într-un anumit tip de canal de comunicație, pachetele de informație și se înregistrează lungimea intervalului dintre două pachete succesive. Numele de populație și indivizi vine din biologie și demografie, domenii în care s-au făcut pentru prima dată investigații statistice.

Definiția 1.1.1 Fie (X, f_θ) un model statistic asociat unei populații. Un vector aleator $\xi = (X_1, X_2, \dots, X_n)$, ale cărui coordonate sunt independente și identic distribuite după legea modelului f , se numește *selecție aleatoare* de volum n . În urma investigării prin sondaj a populației, se înregistrează n valori numerice (x_1, x_2, \dots, x_n) , numite *valori de selecție* sau *valori de realizare* a selecției aleatoare ξ .

Fie $\xi = (X_1, X_2, \dots, X_n)$ o selecție aleatoare de volum n asociată modelului statistic (X, f_θ) . O funcție reală continuă de variabile X_1, X_2, \dots, X_n ,

$$Y = T(X_1, X_2, \dots, X_n),$$

este o variabilă aleatoare, numită *statistică*. Dacă (x_1, x_2, \dots, x_n) este o realizare a selecției aleatoare (X_1, X_2, \dots, X_n) , atunci $T(x_1, x_2, \dots, x_n)$ este o realizare a lui Y . Distribuția de probabilitate a statisticii Y se numește *distribuția de selecție* a statisticii. Această distribuție poate fi dedusă pe baza unor rezultate de teoria probabilităților sau poate fi aproximată.

Având valorile de selecție x_1, x_2, \dots, x_n , primele informații ce se extrag din acestea sunt *media de selecție* sau *media experimentală*, notată cu \bar{x} , care este media lor aritmetică:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

și *dispersia de selecție* (*dispersia experimentală*), s^2 , definită prin:

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

Radical din dispersia de selecție, $\sqrt{s^2}$, se notează cu s și se numește *abaterea standard* a eșantionului.

Media de selecție \bar{x} este realizare a statisticii \bar{X} , unde

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

iar dispersia de selecție este o realizare a statisticii

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

Propoziția 1.1.1 Dacă (X_1, X_2, \dots, X_n) este o selecție aleatoare de volum n dintr-o populație modelată statistic de (X, f) , cu $M(X) = m$ și $\sigma^2(X) = \sigma^2$, atunci statistica medie \bar{X} are și ea aceeași valoare medie m , iar dispersia sa este $D^2 = \sigma^2/n$.

Demonstrație: Variabilele X_k au media $m = M(X_k)$ și dispersia $\sigma^2 = \sigma^2(X_k)$. Astfel,

$$M(\bar{X}) = M\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \frac{M(X_1) + \dots + M(X_n)}{n} = \frac{n m}{n} = m,$$

iar

$$\sigma^2(\bar{X}) = \sigma^2\left(\frac{X_1}{n} + \dots + \frac{X_n}{n}\right) = \frac{1}{n^2}\sigma^2(X_1) + \dots + \frac{1}{n^2}\sigma^2(X_n) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}.$$

În calculul dispersiei am aplicat formula de calcul a dispersiei unei combinații liniare cu coeficienți reali de variabile aleatoare independente:

$$\sigma^2(a_1 X_1 + a_2 X_2 + \dots + a_n X_n) = a_1^2 \sigma^2(X_1) + a_2^2 \sigma^2(X_2) + \dots + a_n^2 \sigma^2(X_n).$$

□

1.2 Estimatori ai parametrilor modelelor statistice

Fie (X, f_θ) un model statistic și (x_1, x_2, \dots, x_n) o realizare a unei selecții aleatoare de volum n , (X_1, X_2, \dots, X_n) .

Un estimator punctual al parametrului θ este o funcție $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$. Este evident că un estimator punctual este o realizare a variabilei aleatoare $\hat{\theta}(X_1, X_2, \dots, X_n)$. Deoarece există o infinitate de funcții $\hat{\theta}$, înseamnă că există o infinitate de estimatori ai parametrului θ . Este însă rezonabil să alegem estimatori care să aproximeze parametrul distribuției f_θ cu o probabilitate suficient de mare.

Definiția 1.2.1 Estimatorul $\hat{\theta}(x_1, x_2, \dots, x_n)$ cu proprietatea că pentru orice $\varepsilon > 0$ are loc

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}(X_1, \dots, X_n) - \theta| > \varepsilon) = 0 \quad (1)$$

se numește *estimator consistent* al parametrului θ .

Intuitiv, dacă estimatorul este consistent, atunci distribuția statisticii $\hat{\theta}(X_1, \dots, X_n)$ este din ce în ce mai concentrată în jurul parametrului θ pe măsură ce volumul selecției crește.

Definiția 1.2.2 Un estimator $\hat{\theta}(x_1, x_2, \dots, x_n)$ care verifică proprietatea că valoarea medie a statisticii $\hat{\theta}(X_1, X_2, \dots, X_n)$ este chiar parametrul θ , adică

$$M(\hat{\theta}(X_1, X_2, \dots, X_n)) = \theta, \quad (2)$$

se numește *estimator centrat* sau *nedeplasat*.

Definiția 1.2.3 Fie $\hat{\theta}_1, \hat{\theta}_2$ doi estimatori nedeplasați ai parametrului θ . Dacă între dispersiile statisticilor $\hat{\theta}_1(X_1, \dots, X_n)$ și $\hat{\theta}_2(X_1, \dots, X_n)$ există relația

$$\sigma^2(\hat{\theta}_1(X_1, \dots, X_n)) \leq \sigma^2(\hat{\theta}_2(X_1, \dots, X_n)), \quad (3)$$

atunci estimatorul $\hat{\theta}_1$ se zice că este *mai eficient* decât estimatorul $\hat{\theta}_2$.

Exemplul 1. Fie (x_1, x_2) o realizare a selecției aleatoare (X_1, X_2) , distribuția comună a variabilelor aleatoare $X_i, i = 1, 2$, fiind distribuția exponențială de parametru θ ,

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

Considerăm trei estimatori pentru θ :

$$\hat{\theta}_1 = x_1, \quad \hat{\theta}_2 = \frac{x_1 + x_2}{2}, \quad \hat{\theta}_3 = \frac{x_1 + x_2}{3}.$$

Estimatorii dați sunt realizări ale variabilelor aleatoare:

$$\begin{aligned} \hat{\theta}_1(X_1, X_2) &= X_1, \\ \hat{\theta}_2(X_1, X_2) &= (X_1 + X_2)/2, \\ \hat{\theta}_3(X_1, X_2) &= (X_1 + X_2)/3. \end{aligned}$$

Se știe că media unei variabile aleatoare exponențial distribuite este θ . Prin urmare, $M(X_1) = \theta$, $M((X_1 + X_2)/2) = (\theta + \theta)/2 = \theta$, iar $M((X_1 + X_2)/3) = 2\theta/3$. Astfel, primii doi estimatori sunt nedeplasați, iar al treilea este deplasat. Să determinăm care este mai eficient dintre primii doi:

$$\sigma^2(X_1) = \theta^2, \quad \sigma^2((X_1 + X_2)/2) = \sigma^2(X_1)/4 + \sigma^2(X_2)/4 = \theta^2/4 + \theta^2/4 = \theta^2/2.$$

În concluzie al doilea estimator este mai eficient.

Existând mai multe posibilități de a defini estimatori ai parametrului, ne întrebăm dacă există o metodă ce permite definirea unui estimator "bun". Răspunsul este pozitiv în câteva cazuri de interes.

1.2.1 Estimarea mediei

Fie (X, f) un model statistic continuu sau discret și x_1, x_2, \dots, x_n observații independente din legea f .

Propoziția 1.2.1 Media de selecție a observațiilor x_1, x_2, \dots, x_n ,

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}, \quad (4)$$

este un estimator nedeplasat al mediei $m = M(X)$ a modelului statistic.

Demonstrație: Rezultă imediat din faptul că $M(\overline{X}) = m$, deci $\hat{m} = \overline{x}$ este un estimator nedeplasat al mediei m . \square

În concluzie, un bun estimator al mediei $M(X)$, a distribuției oricărui model statistic, este media aritmetică a valorilor de selecție.

Exemplul 2. 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 să se estimeze θ în funcție de media de selecție \overline{x} a unei selecții aleatoare de volum n .

b) Să se determine un estimator 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 \overline{x} este media de selecție a unui eșantion de valori x_1, x_2, \dots, x_n , atunci din egalitatea impusă $\hat{m} = \overline{x}$ se determină un estimator al parametrului θ :

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

b) Pentru valorile înregistrate avem $\overline{x} = 0.6$, deci un estimator pentru parametrul θ este $\hat{\theta} = \frac{2\overline{x} - 1}{1 - \overline{x}} = 0.5$.

Exemplul 3. Pentru a estima rata sosirii λ a cererilor de acces la o bază de date s-au monitorizat intervalele de timp dintre 10 cereri consecutive și s-au înregistrat valorile:

$$0.2, 0.1, 0.1, 0.05, 0.05, 0.2, 0.8, 0.5, 0.2, 0.8.$$

Care este estimatorul ratei sosirilor, $\hat{\lambda}$?

Sosirile cererilor la o bază de date este un proces Poisson (N_t) , $t \geq 0$, de rată $\lambda > 0$. Intervalul inter-sosirilor are distribuția exponențială, $X \sim \text{Exp}(\theta = 1/\lambda)$. Dar parametrul θ pentru distribuția exponențială este valoarea medie a variabilei X .

Prin urmare din datele înregistrate putem calcula un estimator al lui θ , adică a mediei lungimii intervalelor inter-sosirilor: $\hat{\theta} = 0.3$. Din relația $\theta = 1/\lambda$, adică $\lambda = 1/\theta$, obținem un estimator al ratei sosirilor ca fiind $\hat{\lambda} = 1/\hat{\theta} = 1/0.3 = 3.33$.

1.2.2 Estimarea dispersiei

Propoziția 1.2.2 Dacă (X, f) este un model statistic și m, σ^2 sunt media și dispersia variabilei aleatoare X , atunci dispersia valorilor de selecție x_1, x_2, \dots, x_n din legea de probabilitate definită de f ,

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad (5)$$

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

Demonstrație: Statistica S^2 , a cărei realizare este s^2 , este dată de relația:

$$S^2 = \hat{\theta}(X_1, X_2, \dots, X_n) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \quad (6)$$

unde $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$. Media acestei statistici este

$$M(S^2) = \frac{1}{n-1} M \left(\sum_{i=1}^n (X_i - \bar{X})^2 \right). \quad (7)$$

Să explicităm $\sum_{i=1}^n (X_i - \bar{X})^2$:

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n [(X_i - m) - (\bar{X} - m)]^2 \\ &= \sum_{i=1}^n ((X_i - m)^2 - 2(\bar{X} - m)(X_i - m) + (\bar{X} - m)^2) \\ &= \sum_{i=1}^n (X_i - m)^2 - 2(\bar{X} - m) \sum_{i=1}^n (X_i - m) + n(\bar{X} - m)^2 \\ &= \sum_{i=1}^n (X_i - m)^2 - 2n(\bar{X} - m)^2 + n(\bar{X} - m)^2 \\ &= \sum_{i=1}^n (X_i - m)^2 - n(\bar{X} - m)^2, \end{aligned} \quad (8)$$

ceea ce implică

$$M \left(\sum_{i=1}^n (X_i - \bar{X})^2 \right) = M \left(\sum_{i=1}^n (X_i - m)^2 \right) - nM((\bar{X} - m)^2). \quad (9)$$

Variabilele aleatoare X_1, X_2, \dots, X_n sunt identic distribuite și, prin urmare, au aceeași medie m și dispersie σ^2 . Deci, $M((X_i - m)^2) = \sigma^2, \forall i = \overline{1, n}$, iar

$$M \left(\sum_{i=1}^n (X_i - m)^2 \right) = \sum_{i=1}^n M((X_i - m)^2) = n\sigma^2.$$

Astfel,

$$M\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) = n\sigma^2 - nD^2(\bar{X}) = n\sigma^2 - n\frac{\sigma^2}{n} = (n-1)\sigma^2. \quad (10)$$

Împărțind la $n-1$, obținem rezultatul dorit. \square

Propoziția 1.2.2 afirmă că media statisticii S^2 este $M(S^2) = \sigma^2$, unde σ^2 este dispersia legii de probabilitate a modelului statistic.

1.3 Estimatorul verosimilității maxime

Estimatorii nedeplasați pentru medie și dispersie s-au determinat simplu, din valorile de selecție dintr-o lege de probabilitate total necunoscută, prin metode ce nu sunt explicit formulate. În cazul în care legea de probabilitate a modelului statistic este cunoscută, în sensul că se cunoaște densitatea de probabilitate, dar aceasta depinde de unul sau mai mulți parametri necunoscuți, estimatorii parametrilor pot fi determinați prin metode bine fundamentate științific. În acest scop considerăm că se supune investigației statistice o caracteristică a unei populații, măsurată de o variabilă aleatoare X a cărei densitate de probabilitate f_θ depinde de un parametru necunoscut θ . Se investighează un eșantion de volum n din populația respectivă și se înregistrează valorile (de observație asupra lui X) x_1, x_2, \dots, x_n . Apoi se determină un estimator pentru parametrul θ , care maximizează probabilitatea înregistrării unor valori foarte apropiate de acestea.

Mai precis, pentru fiecare parametru θ , probabilitatea ca variabila aleatoare X să ia valori apropiate de x_i este probabilitatea ca X să ia valori într-un interval de forma $[x_i, x_i + h)$, cu h foarte mic. Această probabilitate este

$$P(X \in [x_i, x_i + h)) = \int_{x_i}^{x_i+h} f_\theta(x) dx \approx f_\theta(x_i)h.$$

Cu alte cuvinte, aria trapezului curbiliniu de baza h (valoarea integralei) se poate aproxima cu aria dreptunghiului de bază h și înălțime $f_\theta(x_i)$.

Notând cu X_1, X_2, \dots, X_n variabilele aleatoare independente și identic distribuite ca X , avem că probabilitatea ca variabilele X_1, X_2, \dots, X_n să ia valori foarte apropiate de valorile înregistrate x_1, x_2, \dots, x_n este:

$$\begin{aligned} P(X_1 \in [x_1, x_1 + h), X_2 \in [x_2, x_2 + h), \dots, X_n \in [x_n, x_n + h)) \\ = P(X_1 \in [x_1, x_1 + h))P(X_2 \in [x_2, x_2 + h)) \cdots P(X_n \in [x_n, x_n + h)) \\ = f_\theta(x_1)f_\theta(x_2) \cdots f_\theta(x_n)h^n. \end{aligned}$$

Deoarece h^n nu depinde de θ , rezultă că parametrul θ ce maximizează probabilitatea

$$P(X_1 \in [x_1, x_1 + h), X_2 \in [x_2, x_2 + h), \dots, X_n \in [x_n, x_n + h))$$

este parametrul ce maximizează produsul

$$f_\theta(x_1)f_\theta(x_2) \cdots f_\theta(x_n).$$

Definiția 1.3.1 Funcția $L : \mathbb{R} \rightarrow \mathbb{R}$, de variabilă θ , asociată eșantionului x_1, x_2, \dots, x_n , definită prin

$$L(\theta; x_1, x_2, \dots, x_n) = f_\theta(x_1) \cdot f_\theta(x_2) \cdots f_\theta(x_n),$$

se numește *funcția de verosimilitate*.

Valorile x_1, x_2, \dots, x_n fiind cunoscute (fiind valorile rezultate din observații), funcția de verosimilitate este o funcție de o singură variabilă, θ .

Dacă eșantionul s-a extras dintr-o populație a cărei distribuție de probabilitate este discretă și depinde de un parametru θ , adică $p_X(x; \theta) = P(X = x)$, atunci definim funcția de verosimilitate astfel

$$L(\theta; x_1, x_2, \dots, x_n) = p_X(x_1; \theta) p_X(x_2; \theta) \cdots p_X(x_n; \theta).$$

În acest caz, $L(\theta; x_1, x_2, \dots, x_n)$ reprezintă, datorită independenței variabilelor discrete X_1, X_2, \dots, X_n , probabilitatea ca X_1 să ia valoarea x_1 , X_2 să ia valoarea x_2, \dots, X_n să ia valoarea x_n .

Atât în cazul continuu, cât și în cel discret, **estimatorul verosimilității maxime** a parametrului θ este

$$\hat{\theta} = \operatorname{argmax}(L(\theta; x_1, x_2, \dots, x_n)),$$

unde prin $\operatorname{argmax}(L(\theta; x_1, x_2, \dots, x_n))$ se înțelege argumentul θ care maximizează funcția L .

Exemplul 4. Fie populația \mathcal{P} formată dintr-un tip de circuite. Caracteristica ce dorim să o investigăm prin sondaj statistic este durata de viață a acestor circuite, știind că aceasta este exponențial distribuită, cu parametrul θ necunoscut. Măsurând timpul de viață (în ani) a 10 circuite, se obțin valorile:

$$0.8830, 1.96511, 1.9189, 4.8448, 0.9208, 3.4377, 1.7162, 4.2327, 5.9435, 8.3128.$$

Să determinăm estimatorul de verosimilitate maximă pentru θ (adică pentru media duratei de viață a acestui tip de circuite). Densitatea de probabilitate a distribuției exponențiale este

$$f_\theta(x) = \begin{cases} 0, & \text{dacă } x < 0, \\ \frac{1}{\theta} e^{-x/\theta}, & \text{dacă } x \geq 0. \end{cases} \quad (11)$$

Astfel, funcția de verosimilitate este

$$L(\theta; x_1, x_2, \dots, x_n) = \prod_{i=1}^n \frac{1}{\theta} e^{-x_i/\theta} = \frac{1}{\theta^n} e^{-\sum_{i=1}^n x_i/\theta}. \quad (12)$$

Funcția logaritmică cu bază mai mare ca 1 are derivata întâi pozitivă. Notând cu $h = \ln$ și cu $\ell(\theta) = h(L(\theta))$, avem că funcția $\ell'(\theta) = h(L(\theta))L'(\theta)$ are același semn ca derivata lui L , deci ℓ și L au aceleași puncte de extrem și de aceeași natură. Pentru simplitatea calculelor, determinăm punctul de maxim absolut (dacă acesta există) pentru ℓ și acesta va fi punct de maxim absolut și pentru L :

$$\ell(\theta) = \ln L(\theta; x_1, x_2, \dots, x_n) = -n \ln \theta - \frac{\sum_{i=1}^n x_i}{\theta}. \quad (13)$$

Avem

$$\ell'(\theta) = -\frac{n}{\theta} + \frac{\sum_{i=1}^n x_i}{\theta^2}. \quad (14)$$

Rezolvând ecuația $\ell'(\theta) = 0$ în raport cu θ , obținem punctul $\hat{\theta} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$, care este maxim absolut pentru ℓ , deci și pentru L . În concluzie,

$$\operatorname{argmax} L(\theta; x_1, x_2, \dots, x_n) = \bar{x},$$

estimatorul de verosimilitate maximă a parametrului θ a distribuției exponențiale este media de selecție. În cazul exemplului dat, estimatorul verosimilității maxime a mediei de viață a circuitelor este media selecției:

$$\hat{\theta} = \frac{x_1 + x_2 + \dots + x_{10}}{10} = 5.1861.$$

Observație: Valorile de selecție au fost generate simulând o variabilă $X \sim \operatorname{Exp}(\theta = 5)$, deci estimatorul verosimilității maxime $\hat{\theta} = 5.1861$ este ”destul de bun”.

Exemplul 5. Un simulator al distribuției Bernoulli

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

de parametru p necunoscut, generează stringul de biți:

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

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

Rezolvare: Notând parametrul p necunoscut cu θ , distribuția de probabilitate a variabilei X este $p_X(\theta; b) = P(X = b)$, unde b este bitul 1 sau 0.

$$p_X(\theta; b) = \begin{cases} \theta, & \text{dacă } b = 1, \\ 1 - \theta, & \text{dacă } b = 0. \end{cases}$$

Funcția de verosimilitate asociată eșantionului de biți generați este

$$L(\theta; b_1, b_2, \dots, b_{26}) = p_X(\theta; b_1) p_X(\theta; b_2) \cdots p_X(\theta; b_{26}) = \theta^{12} (1 - \theta)^{14},$$

unde 12 este numărul de biți 1 din string, iar 14 numărul de biți 0. Se logaritmează și se determină punctul de maxim absolut al funcției $\ell(\theta) = \ln(L(\theta))$. Cum

$$\ell(\theta) = 12 \ln(\theta) + 14 \ln(1 - \theta),$$

rezultă

$$\ell'(\theta) = \frac{12}{\theta} - \frac{14}{1 - \theta} = \frac{12(1 - \theta) - 14\theta}{\theta(1 - \theta)}.$$

Ecuția $\ell'(\theta) = 0$ are soluția $\hat{\theta} = \frac{12}{26}$. Se verifică că acesta este un punct de maxim pentru ℓ , deci estimatorul verosimilității maxime pentru parametrul p al distribuției Bernoulli, dedus din stringul de biți, este egal cu numărul biților 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 din string supra numărul cazurilor posibile.

1.3.1 Estimarea mediei și dispersiei distribuției normale

Distribuția normală apare în numeroase modele statistice și estimarea parametrilor ei intervine, de exemplu, în numeroase probleme de *machine learning*. Tocmai de aceea, prezintă interes determinarea estimatorului de verosimilitate maximă atât pentru media, cât și pentru dispersia lui $N(m, \sigma^2)$. Discutăm trei cazuri:

1. Fie $(X, f(x; m, \sigma^2))$ un model statistic caracterizat de distribuția normală de **medie** m **necunoscută** și **dispersie cunoscută**. Densitatea de probabilitate a distribuției normale este

$$f(x; m) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-m)^2/(2\sigma^2)}, \quad x \in \mathbb{R}. \quad (15)$$

Funcția de verosimilitate asociată realizării (x_1, x_2, \dots, x_n) a unei selecții aleatoare de volum n este

$$L(m) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - m)^2}{2\sigma^2}} = C e^{-\sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}}, \quad m \in \mathbb{R}, \quad (16)$$

unde C este o constantă pozitivă ce nu depinde de m . Pentru a determina punctele de extrem ale lui L considerăm funcția ℓ , $\ell = \ln(L)$:

$$\ell(m) = \ln C - \sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}. \quad (17)$$

Calculăm

$$\ell'(m) = \frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - m) = \frac{1}{\sigma^2} (x_1 + x_2 + \dots + x_n - nm). \quad (18)$$

Punctul $\hat{m}(x_1, x_2, \dots, x_n) = (x_1 + x_2 + \dots + x_n)/n$ este punctul de maxim pentru ℓ , deci și pentru L . Prin urmare, *media selecției este estimatorul de verosimilitate maximă pentru media distribuției normale*.

2. În cazul în care media este cunoscută și dispersia necunoscută, spațiul parametrului σ^2 este $\Theta = (0, \infty)$. Printr-un calcul analog, rezultă că estimatorul de verosimilitate maximă al dispersiei este

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - m)^2. \quad (19)$$

Evident că acest estimator este deplasat, adică

$$M \left(\frac{1}{n} \sum_{i=1}^n (X_i - m)^2 \right) \neq \sigma^2. \quad (20)$$

Un estimator nedeplasat pentru dispersie este

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (21)$$

Se spune că s^2 s-a obținut din $\hat{\sigma}^2$ prin ajustare.

3. A treia alternativă este atunci când atât media, cât și dispersia sunt necunoscute. În acest caz funcția de verosimilitate depinde de doi parametri:

$$L(m, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}}. \quad (22)$$

Logaritmând din nou, obținem

$$\ell(m, \sigma^2) = -\frac{n}{2} \ln(2\pi\sigma^2) - \sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}. \quad (23)$$

Punctele staționare ale lui ℓ sunt soluții ale sistemului:

$$\frac{\partial \ell}{\partial m} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - m) = 0 \quad (24)$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - m)^2 = 0, \quad (25)$$

adică

$$\begin{aligned} \hat{m}(x_1, x_2, \dots, x_n) &= \frac{x_1 + x_2 + \dots + x_n}{n} = \bar{x}, \\ \hat{\sigma}^2(x_1, x_2, \dots, x_n) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2. \end{aligned}$$

Se verifică apoi că $(\hat{m}, \hat{\sigma}^2)$ este punct de maxim pentru ℓ , deci și pentru L , arătând că $\frac{\partial^2 \ell}{\partial m^2}(\hat{m}, \hat{\sigma}^2) < 0$ și

$$\begin{vmatrix} \frac{\partial^2 \ell}{\partial m^2} & \frac{\partial^2 \ell}{\partial m \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial m \partial \sigma^2} & \frac{\partial^2 \ell}{\partial (\sigma^2)^2} \end{vmatrix}(\hat{m}, \hat{\sigma}^2) > 0. \quad (26)$$

Estimatorul de maximă verosimilitate pentru medie este nedeplasat, dar pentru dispersie este deplasat. Ajustând estimatorul $\hat{\sigma}^2$ la

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

obținem estimator nedeplasat.

1.4 Teorema limită centrală

În statistică prezintă importanță deosebită variabila medie aritmetică (media de selecție), notată \bar{X} , asociată unei selecții aleatoare (X_1, X_2, \dots, X_n) ,

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}. \quad (27)$$

Dacă x_1, x_2, \dots, x_n sunt valori de selecție, atunci media de selecție \bar{x} este o observație asupra variabilei \bar{X} . Variabila aleatoare medie aritmetică \bar{X} are media $M(\bar{X}) = m$ și dispersia $D^2(\bar{X}) = \sigma^2/n$. Remarcăm că dispersia mediei aritmetice descrește invers proporțional cu n , iar abaterea standard descrește invers proporțional cu \sqrt{n} . Cum abaterea standard este o măsură a împrăstierii valorilor lui \bar{X} , relația $D(\bar{X}) = \sigma/\sqrt{n}$ arată că pe măsură ce n crește, distribuția mediei aritmetice este din ce în ce mai concentrată în jurul valorii medii m . De exemplu, dacă $n = 100$, $m = 0$ și $\sigma = 3$, atunci $D(\bar{X}) = 3/10 = 0.3$. Crescând n la 1000, $D(\bar{X}) = 3/\sqrt{1000} = 0.0949$.

Propoziția 1.4.1 *Dacă X_1, X_2, \dots, X_n sunt variabile aleatoare independente și normal distribuite, $X_i \sim N(m_i, \sigma_i^2)$, atunci pentru orice $a_i \in \mathbb{R}$, $i = \overline{1, n}$, combinația liniară a variabilelor aleatoare cu coeficienții a_i este normal distribuită,*

$$X = a_1 X_1 + a_2 X_2 + \dots + a_n X_n \sim N \left(\sum_{i=1}^n a_i m_i, \sum_{i=1}^n a_i^2 \sigma_i^2 \right). \quad (28)$$

Precizare: Media și dispersia oricărei combinații liniare de variabile independente se calculează la fel:

$$M(a_1 X_1 + \dots + a_n X_n) = a_1 M(X_1) + \dots + a_n M(X_n),$$

$$\sigma^2(a_1X_1 + \dots + a_nX_n) = a_1^2\sigma^2(X_1) + \dots + a_n^2\sigma^2(X_n).$$

Informația adițională pe care o aduce Propoziția 1.4.1 este că dacă variabilele aleatoare ce intră în combinație au distribuția normală, atunci și combinația lor liniară are distribuția normală.

Exemplul 6. Unui semnal X , transmis printr-un canal de comunicație, i se adaugă un zgomot N . Știind că variabilele aleatoare X și N sunt independente și $X \sim N(0, 1)$, $N \sim N(0, \sigma^2)$, să se determine distribuția de probabilitate a semnalului $Y = X + N$ detectat.

Rezolvare: Y fiind o combinație liniară cu coeficienții $a_1 = a_2 = 1$ a variabilelor aleatoare independente și normal distribuite X și N , rezultă că $Y \sim N(0, 1 + \sigma^2)$, unde $1 + \sigma^2$ reprezintă dispersia variabilei Y .

În statistică este important următorul rezultat:

Corolar 1.4.1 *Dacă variabilele aleatoare i.i.d. X_1, X_2, \dots, X_n au distribuția normală, $X_i \sim N(m, \sigma^2)$, $i = \overline{1, n}$, atunci media lor aritmetică are de asemenea distribuția normală, cu aceeași medie m și dispersie $\frac{\sigma^2}{n}$, adică*

$$\overline{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(m, \sigma^2/n). \quad (29)$$

Ne întrebăm în mod natural ce distribuție de probabilitate are media aritmetică a n variabile aleatoare i.i.d. având distribuția comună de probabilitate absolut arbitrară, adică ne-normală (ea putând fi exponențială, binomială etc). Răspunsul este dat de unul din cele mai remarcabile rezultate din teoria probabilităților, cu aplicații importante în statistică:

Teorema 1.4.1 (*Teorema limită centrală*) *Dacă (X_n) este un șir de variabile aleatoare independente și identic distribuite având media comună m și abaterea standard σ , iar*

$$\overline{X}_n = \frac{1}{n}(X_1 + X_2 + \dots + X_n) \quad (30)$$

este șirul variabilelor medie aritmetică, atunci pentru $n \rightarrow \infty$ distribuția de probabilitate a variabilelor \overline{X}_n este aproximativ normală de medie m și dispersie $D^2 = \sigma^2/n$. Notăm $\overline{X}_n \sim ApN(m, D^2 = \sigma^2/n)$.

Cu alte cuvinte, teorema precedentă afirmă că "în medie totul este normal". Formulată riguros (matematic), teorema limită centrală asigură în condițiile enunțate că funcțiile de repartiție ale variabilelor standardizate

$$Z_n = \frac{\overline{X}_n - m}{\frac{\sigma}{\sqrt{n}}}$$

tind, când $n \rightarrow \infty$, la funcția de repartiție Φ a distribuției normale standard, $N(0, 1)$, adică

$$\lim_{n \rightarrow \infty} P(Z_n \leq x) = \Phi(x).$$

În practica statistică se consideră că pentru $n \geq 30$, distribuția normală poate fi folosită ca distribuție a mediei aritmetice a n variabile aleatoare i.i.d. cu media și dispersia finită. Cu alte cuvinte, dacă x_1, x_2, \dots, x_n este un eșantion de volum $n \geq 30$ dintr-o populație a cărei caracteristică de interes are o distribuție de probabilitate arbitrară, de medie m și abatere standard σ , media de selecție \bar{x} poate fi considerată ca o observație asupra unei variabile aleatoare normal distribuite de medie m și abatere standard σ/\sqrt{n} .

Pentru a verifica această proprietate efectuăm următorul experiment: generăm de 100 de ori consecutiv câte n numere aleatoare din distribuția exponențială de parametru $\theta = 2.3$. Asociem fiecărui șir de numere aleatoare

$$x_1^k, x_2^k, \dots, x_n^k, \quad k = \overline{1, 100}, \quad (31)$$

media aritmetică a valorilor sale

$$\bar{x}^k = \frac{x_1^k + x_2^k + \dots + x_n^k}{n}, \quad k = \overline{1, 100}. \quad (32)$$

Valorile $\bar{x}^1, \bar{x}^2, \dots, \bar{x}^{100}$ le interpretăm ca observații asupra variabilei aleatoare

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

adică a mediei de selecție a n variabile aleatoare i.i.d. $X_1, X_2, \dots, X_n \sim \text{Exp}(\theta = 2.3)$. Generăm apoi funcția de repartiție empirică a datelor $\bar{x}^1, \bar{x}^2, \dots, \bar{x}^{100}$ și afișăm simultan graficul acesteia, comparativ cu graficul repartiției Φ a distribuției normale standard. Se observă că pe măsură ce volumul n al eșantioanelor generate crește, funcția de repartiție empirică este din ce în ce mai bine aproximată de repartiția Φ .

O altă metodă de a vizualiza rezultatul teoremei limită centrală este de a genera histograma datelor $\bar{x}^1, \bar{x}^2, \dots, \bar{x}^{100}$, care pe măsură ce n crește, aproximează clopotul lui Gauss.

Teorema limită centrală prezintă interes și în următorul context: șirului de variabile aleatoare i.i.d. (X_n) îi asociem șirul (S_n) , definit prin

$$S_n = X_1 + X_2 + \dots + X_n.$$

Evident, $S_n = n\bar{X}_n$, $M(S_n) = M(n\bar{X}_n) = nM(\bar{X}_n) = nm$, iar

$$\sigma^2(S_n) = \sigma^2(n\bar{X}_n) = n^2\sigma^2(\bar{X}_n) = n^2\sigma^2/n = n\sigma^2.$$

Prin urmare, pentru n suficient de mare, S_n fiind combinație liniară a unor variabile aleatoare aproximativ normal distribuite, este și ea aproximativ normal distribuită:

$$S_n \sim \text{ApN}(nm, D^2 = n\sigma^2). \quad (33)$$

Exemplul 7. Variabila aleatoare X , care dă numărul de zone defecte ale unui CD de un anumit tip, are următoarea distribuție de probabilitate: $p(x) = P(X = x)$, unde

x	$p(x)$
0	0.75
1	0.15
2	0.10

- a) Să se calculeze media și abaterea standard a numărului de zone defecte ale CD-ului.
b) Ce distribuție de probabilitate are media de selecție a unui eșantion de 400 de CD-uri din tipul investigat? Care este media și dispersia acestei distribuții?
c) Care este probabilitatea ca media numărului de zone defecte/CD într-un lot de 400 de CD-uri să fie mai mică decât 0.3?

Rezolvare: a) $m = M(X) = 0.15 + 0.2 = 0.35$, iar

$$\sigma^2(X) = (0 - 0.35)^2 0.75 + (1 - 0.35)^2 0.15 + (2 - 0.35)^2 0.1 = 0.4275$$

și deci abaterea standard este $\sigma = \sigma(X) = \sqrt{\sigma^2(X)} = 0.6538$.

b) Volumul eșantionului fiind mare, conform teoremei limită centrală

$$\overline{X}_{400} \sim \text{ApN}(m, \sigma^2/400),$$

unde $m = 0.35$, $D^2 = \sigma^2/400 = 0.00106$, iar $D = 0.032$.

c) Avem

$$\begin{aligned} P(\overline{X}_{400} < 0.3) &= F_{\overline{X}}(0.3) = \Phi\left(\frac{0.3 - m}{D}\right) = \Phi\left(\frac{0.3 - 0.35}{0.032}\right) \\ &= \Phi(-1.5625) = 1 - \Phi(1.5625) = 1 - 0.94 = 0.06. \end{aligned}$$

Exemplul 8. Se consideră o buclă `for`:

```
for(i=1; i<=n; i++) // n>30
{
    executa blocul B;
}
```

Știind că timpul de execuție al blocului B este o variabilă aleatoare de distribuție de probabilitate necunoscută, având media $m = 60ms$ și abaterea standard de $\sigma = 8ms$, iar execuțiile succesive ale blocului sunt independente, să se determine distribuția de probabilitate a timpului de execuție a buclei `for`. Care este probabilitatea ca timpul de execuție al buclei în cazul $n = 50$ să fie cuprins între $0.75s$ sec și 1 sec?

Notând cu T_i timpul celei de-a i -a execuții a blocului B, timpul total de execuție

$$T = T_1 + T_2 + \cdots + T_{50}$$

este aproximativ normal distribuit.