

# Cursul 11

## Introducere în Statistică

### 11.1 Problematica 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, 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  $f_\theta$  a variabilei aleatoare continue  $X$  sau distribuția de probabilitate  $p_\theta$ , 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  $f_\theta$  sau  $p_\theta$ .

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  $(X, f_\theta)$ , unde  $X$  este o variabilă aleatoare reală de densitate  $f_\theta$ , 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  $(X, f_\theta)$  cu  $X \sim \text{Exp}(\theta)$ ,  $\theta$  fiind un parametru necunoscut. Înregistrând durata de viață a  $n$  chipuri selectate la întâmplare din producția dintr-o anumită perioadă, se va estima parametrul  $\theta$ , care se știe că reprezintă media variabilei aleatoare  $X$  a modelului exponențial. Având un estimator al lui  $\theta$ , 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 intervalelor 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 11.1.1** Fie  $(X, f_\theta)$  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  $\xi$ .

Fie  $\xi = (X_1, X_2, \dots, X_n)$  o selecție aleatoare de volum  $n$  asociată modelului statistic  $(X, f_\theta)$ . 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 medie aritmetică  $\bar{X}$ , 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 11.1.1** 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 aritmetică  $\bar{X}$  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,

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

iar

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

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

$$\sigma^2(a_1X_1 + a_2X_2 + \cdots + a_nX_n) = a_1^2\sigma^2(X_1) + a_2^2\sigma^2(X_2) + \cdots + a_n^2\sigma^2(X_n).$$

□

## 11.2 Teorema limită centrală

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

$$\bar{X} = \frac{X_1 + X_2 + \cdots + X_n}{n}. \quad (11.1)$$

Am arătat mai sus că variabila aleatoare medie aritmetică  $\bar{X}$  are media  $M(\bar{X}) = m$  și dispersia  $\sigma^2(\bar{X}) = \sigma^2/n$ , unde  $M(X_k) = m$  și  $\sigma^2(X_k) = \sigma^2$ ,  $k = \overline{1, n}$ . 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  $\sigma(\bar{X}) = \sigma/\sqrt{n}$  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  $\sigma(\bar{X}) = 3/10 = 0.3$ . Crescând  $n$  la 1000, avem  $\sigma(\bar{X}) = 3/\sqrt{1000} = 0.0949$ .

**Propoziția 11.2.1** *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,*

$$X = a_1X_1 + a_2X_2 + \cdots + a_nX_n \sim N\left(m = \sum_{i=1}^n a_i m_i, \sigma^2 = \sum_{i=1}^n a_i^2 \sigma_i^2\right). \quad (11.2)$$

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

**Exemplul 1.** 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 11.2.1** *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*

$$\overline{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(m, \sigma^2/n). \quad (11.3)$$

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

$$\overline{X}_n = \frac{1}{n}(X_1 + X_2 + \dots + X_n) \quad (11.4)$$

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

Teorema limită centrală afirmă de fapt că ”în medie totul este normal”. Formulată riguros (matematic), aceasta asigură (în condițiile enunțate) că șirul funcțiilor de repartiție ale variabilelor standardizate,

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

ține, când  $n \rightarrow \infty$ , la funcția de repartiție  $\Phi$  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  $\sigma$ , media de selecție  $\bar{x}$  poate fi considerată ca o observație asupra unei variabile aleatoare normal distribuite de medie  $m$  și abatere standard  $\sigma/\sqrt{n}$ .

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,

$$S_n \sim ApN(nm, D^2 = n\sigma^2). \quad (11.5)$$

### 11.3 Estimatori ai parametrilor modelelor statistice

Fie  $(X, f_\theta)$  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  $\theta$  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  $\theta$ . Este însă rezonabil să alegem estimatori care să aproximeze parametrul distribuției  $f_\theta$  cu o probabilitate suficient de mare.

**Definiția 11.3.1** Estimatorul  $\hat{\theta}(x_1, x_2, \dots, x_n)$  cu proprietatea că pentru orice  $\varepsilon > 0$  are loc

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}(X_1, \dots, X_n) - \theta| > \varepsilon) = 0 \quad (11.6)$$

se numește *estimator consistent* al parametrului  $\theta$ .

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  $\theta$  pe măsură ce volumul selecției crește.

**Definiția 11.3.2** 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  $\theta$ , adică

$$M(\hat{\theta}(X_1, X_2, \dots, X_n)) = \theta, \quad (11.7)$$

se numește *estimator centrat* sau *nedeplasat*.

**Definiția 11.3.3** Fie  $\hat{\theta}_1, \hat{\theta}_2$  doi estimatori nedeplasați ai parametrului  $\theta$ . Dacă între dispersiile statisticilor  $\hat{\theta}_1(X_1, \dots, X_n)$  și  $\hat{\theta}_2(X_1, \dots, X_n)$  există relația

$$\sigma^2(\hat{\theta}_1(X_1, \dots, X_n)) \leq \sigma^2(\hat{\theta}_2(X_1, \dots, X_n)), \quad (11.8)$$

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

**Exemplul 2.** 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  $\theta$ ,

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

Considerăm trei estimatori pentru  $\theta$ :

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

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  $\theta$ . 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:

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

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.

### 11.3.1 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 11.3.1** Media de selecție a observațiilor  $x_1, x_2, \dots, x_n$ ,

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}, \quad (11.9)$$

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

**Demonstrație:** Rezultă imediat din faptul că  $M(\bar{X}) = m$ , deci  $\hat{m} = \bar{x}$  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.

**11.3.2 Estimarea dispersiei**

**Propoziția 11.3.2** Dacă  $(X, f)$  este un model statistic și  $m, \sigma^2$  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  $f$ ,

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad (11.10)$$

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

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

$$S^2 = \hat{\theta}(X_1, X_2, \dots, X_n) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \quad (11.11)$$

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

$$M(S^2) = \frac{1}{n-1} M \left( \sum_{i=1}^n (X_i - \bar{X})^2 \right). \quad (11.12)$$

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

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n [(X_i - m) - (\bar{X} - m)]^2 \\ &= \sum_{i=1}^n ((X_i - m)^2 - 2(\bar{X} - m)(X_i - m) + (\bar{X} - m)^2) \\ &= \sum_{i=1}^n (X_i - m)^2 - 2(\bar{X} - m) \sum_{i=1}^n (X_i - m) + n(\bar{X} - m)^2 \\ &= \sum_{i=1}^n (X_i - m)^2 - 2n(\bar{X} - m)^2 + n(\bar{X} - m)^2 \\ &= \sum_{i=1}^n (X_i - m)^2 - n(\bar{X} - m)^2, \end{aligned} \quad (11.13)$$

ceea ce implică

$$M \left( \sum_{i=1}^n (X_i - \bar{X})^2 \right) = M \left( \sum_{i=1}^n (X_i - m)^2 \right) - nM((\bar{X} - m)^2). \quad (11.14)$$

Variabilele aleatoare  $X_1, X_2, \dots, X_n$  sunt identic distribuite și, prin urmare, au aceeași medie  $m$  și dispersie  $\sigma^2$ . 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,

$$M\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) = n\sigma^2 - nD^2(\bar{X}) = n\sigma^2 - n\frac{\sigma^2}{n} = (n-1)\sigma^2. \quad (11.15)$$

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

Propoziția 11.3.2 afirmă că media statisticii  $S^2$  este  $M(S^2) = \sigma^2$ , unde  $\sigma^2$  este dispersia legii de probabilitate a modelului statistic.

## 11.4 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  $f_\theta$  depinde de un parametru necunoscut  $\theta \in \Theta$ . 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  $\theta$ , care maximizează probabilitatea înregistrării unor valori foarte apropiate de acestea.

Mai precis, pentru fiecare parametru  $\theta$ , 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 bază  $h$  (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  $\theta$ , rezultă că parametrul  $\theta$  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 11.4.1** Funcția  $L : \Theta \rightarrow \mathbb{R}$ , de variabilă  $\theta$ , 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 variabilă  $\theta$ .

Dacă eșantionul s-a extras dintr-o populație a cărei distribuție de probabilitate este discretă și depinde de un parametru  $\theta$ , 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  $\theta$  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  $\theta$  care maximizează funcția  $L$  (dacă acesta există).

#### 11.4.1 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

$$f(x; m) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-m)^2/(2\sigma^2)}, \quad x \in \mathbb{R}. \quad (11.16)$$

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

$$L(m) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - m)^2}{2\sigma^2}} = C e^{-\sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}}, \quad m \in \mathbb{R}, \quad (11.17)$$

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

$$\ell(m) = \ln C - \sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}. \quad (11.18)$$

Calculăm

$$\ell'(m) = \frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - m) = \frac{1}{\sigma^2} (x_1 + x_2 + \cdots + x_n - nm). \quad (11.19)$$

Punctul  $\hat{m}(x_1, x_2, \dots, x_n) = (x_1 + x_2 + \cdots + x_n)/n$  este punctul de maxim pentru  $\ell$ , 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  $\sigma^2$  este  $\Theta = (0, \infty)$ . Printr-un calcul analog, rezultă că estimatorul de verosimilitate maximă al dispersiei este

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - m)^2. \quad (11.20)$$

Evident că acest estimator este deplasat, adică

$$M \left( \frac{1}{n} \sum_{i=1}^n (X_i - m)^2 \right) \neq \sigma^2. \quad (11.21)$$

Un estimator nedeplasat pentru dispersie este

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (11.22)$$

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:

$$L(m, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}}. \quad (11.23)$$

Logaritmând din nou, obținem

$$\ell(m, \sigma^2) = -\frac{n}{2} \ln(2\pi\sigma^2) - \sum_{i=1}^n \frac{(x_i - m)^2}{2\sigma^2}. \quad (11.24)$$

Punctele staționare ale lui  $\ell$  sunt soluții ale sistemului:

$$\frac{\partial \ell}{\partial m} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - m) = 0 \quad (11.25)$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - m)^2 = 0, \quad (11.26)$$

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  $\ell$ , deci și pentru  $L$ , arătând că  $\frac{\partial^2 \ell}{\partial m^2}(\hat{m}, \hat{\sigma}^2) < 0$  și

$$\begin{vmatrix} \frac{\partial^2 \ell}{\partial m^2} & \frac{\partial^2 \ell}{\partial m \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial m \partial \sigma^2} & \frac{\partial^2 \ell}{\partial (\sigma^2)^2} \end{vmatrix}(\hat{m}, \hat{\sigma}^2) > 0. \quad (11.27)$$

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.