

## 5.5. Estimatori statistici și momente

În cea mai mare parte a cazurilor este dificil să se lucreze cu funcția densitate de probabilitate pentru un vector aleator multidimensional. De fapt, într-un număr mare de aplicații practice funcția densitate de probabilitate (sau funcția masă de probabilitate) nici nu poate fi cunoscută. În astfel de situații este posibil însă să se lucreze cu momentele distribuției (în mod uzual cu cele de ordin I și II, acestea fiind cantități definite în strânsă legătură cu funcția densitate de probabilitate, *ale cărei trăsături și proprietăți le reflectă*).

După cum vom vedea, cu toate că aceste momente sunt definite, *teoretic*, în strânsă legătură cu funcția densitate de probabilitate, *în practică* nu numai că nu este necesar să se cunoască această funcție pentru estimarea lor dar, în plus, o estimare a funcției densitate de probabilitate poate fi obținută folosindu-ne de estimările acestor momente precum și de presupunerea inițială privind forma parametrică a densității. Această legătură dintre momentele unei v.a.  $x$  și funcția densitate de probabilitate a lui  $x$  ne permite ca, prin intermediul momentelor de ordin I și II să realizăm diferite analize și să tragem concluzii privind caracteristicile vectorului aleator fără a cunoaște efectiv funcția de distribuție ce-l caracterizează. O tratare a acestor aspecte și a semnificației pe care o au aceste momente o vom face atunci când vom discuta *estimarea parametrică punctuală* (vezi în cadrul **Subcapitolul 5.5.1.** paragraful „**Estimarea parametrică punctuală clasică**”).

Pentru a înțelege ce implică, în general, un **proces de estimare statistică** (văzut ca formă tradițională de inferență statistică<sup>36</sup>), precum și, care sunt metodele aferente, de implementare, de care se dispune la ora actuală, facem în cele ce urmează o prezentare de ansamblu a problemei, urmând ca, ulterior, să abordăm doar parte dintre metodele amintite.

---

<sup>36</sup> Reamintim aici că, prin **inferență statistică** se înțelege obținerea de concluzii bazate pe o evidență statistică (informații derivate dintr-un eșantion). Concluziile care se trag sunt asupra caracteristicilor populației din care provine eșantionul sau, mai general, asupra procesului aleator al cărui comportament a fost observat într-o perioadă finită de timp.

Cea de a doua formă tradițională de inferență statistică este *testarea ipotezelor* (numită și testarea statistică de semnificație) și ea, spre deosebire de estimarea statistică, oferă un răspuns de tipul „da”/„nu” unor întrebări statistice.

### 5.5.1. Procesul estimării statistice

*Probabilitatea statistică* este cea care furnizează *soluții optimale*<sup>37</sup> în cazul datelor măsurate ce sunt afectate de zgomot, caracterizate de incertitudine sau sunt incomplete, ea extrăgând din date cea mai multă informație. În prezent există numeroase domenii în care *teoria estimării* își găsește aplicabilitate. Astfel, *procesarea de semnal, controlul calității, telecomunicațiile, interpretarea științifică a experimentelor, teoria controlului*, etc. – sunt doar câteva dintre exemplele pe care le întâlnim frecvent în practică.

**Estimarea**, ca *ramură a statisticii și a procesării de semnal*, se ocupă cu *estimarea* (aproximarea), pe baza datelor empirice, *a valorilor unor parametri* ai modelelor statistice sau chiar *a structurii modelelor statistice* ce descriu *pattern*-urile dintr-o populație aflată în studiu.

Cunoașterea și specificarea, sau nu, *a priori* a structurii modelelor statistice în procesul de estimare împarte aceste modele statistice (respectiv, metodele de estimare) în:

- *parametrice;*
- *neparametrice;*
- *semiparametrice.*

Pentru a înțelege mai bine distincția dintre acestea (precum și terminologia adoptată – *parametric-neparametric*) plecăm de la definiția unui **model statistic**.

Un **model statistic** constă într-un set de ecuații matematice ce descriu comportamentul unui obiect de studiu în termenii *variabilelor aleatoare* și ai *distribuțiilor de probabilitate* (respectiv, densităților de probabilitate) asociate acestora. În consecință, cunoașterea sau nu a structurii modelului se reduce, practic, la cunoașterea sau nu a formei funcționale a densității de probabilitate ce caracterizează v.a. sau vectorul aleator ce descrie populația. O determinare completă a acestei distribuții de probabilitate implică estimarea parametrilor acesteia pe baza eșantionului de date empirice.

Mai general, un model statistic include însă, pe lângă caracterizarea statistică a datelor numerice (definiția de mai sus), aspecte cum ar fi estimarea comportamentului probabilistic viitor al unui sistem pe baza comportamentului acestuia din trecut, extrapolarea sau interpolarea datelor pe baza unei așa-zise celei mai bune aproximări sau a estimațiilor erorii calculate

---

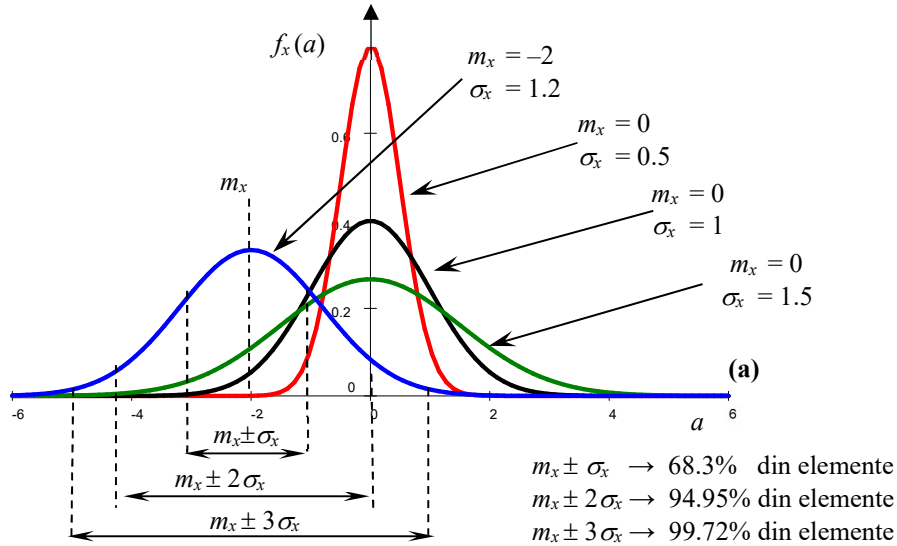
<sup>37</sup> Un *estimator optimal* este estimatorul care extrage din datele măsurate toată informația disponibilă în acestea; în cazul în care în datele empirice mai există informație neutilizată atunci spunem că estimatorul nu este unul optimal.

pentru datele empirice precum și analiza spectrală a datelor unui eșantion statistic.

## 1. Estimarea parametrică

### Estimarea parametrică a densității de probabilitate

În statistică există definite câteva familii de repartiții teoretice remarcabile (împreună cu funcțiile densitate de probabilitate derivate; pentru exemplificare vezi **Anexa: Repartiții remarcabile** și **Subcapitolul 5.6**) iar diferențierea în cadrul unei familii se face printr-un număr finit de parametri (unul sau mai mulți) ce caracterizează respectiva familie de distribuții.



**Figura 5.15.** Exemple de funcții densitate de probabilitate *Gauss*-iene din familia distribuției normale monodimensionale

**Parametrii** în statistică, la fel ca și în matematică, sunt cantități referite prin simboluri ce fac parte din *definiția* funcțiilor și ei definesc anumite caracteristici ale acestora (de exemplu: poziționarea sau alura funcțiilor). Astfel, spre exemplu, în cazul distribuției normale a unei v.a.  $x$ , cei doi parametri ai repartiției,  $m_x$  și  $\sigma_x^2$  (relația (5.115)), au semnificația poziționării distribuției pe axa absciselor (dată de parametrul  $m_x$ ), respectiv, semnificația caracterizării zonelor de concentrare ale punctelor (dată de parametrul  $\sigma_x^2$ ) (vezi **Figura 5.15**, pentru exemplificare cu diverse valori particulare ale celor doi parametri):

$$f_x(a) = f_x(a; m_x, \sigma_x^2) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(a-m_x)^2}{2\sigma_x^2}} \quad (5.115)$$

În relația de mai sus  $a$  este variabila generică a funcției densitate de probabilitate,  $f_x(a)$ , iar  $m_x$  și  $\sigma_x^2$  sunt parametrii acestei funcții.

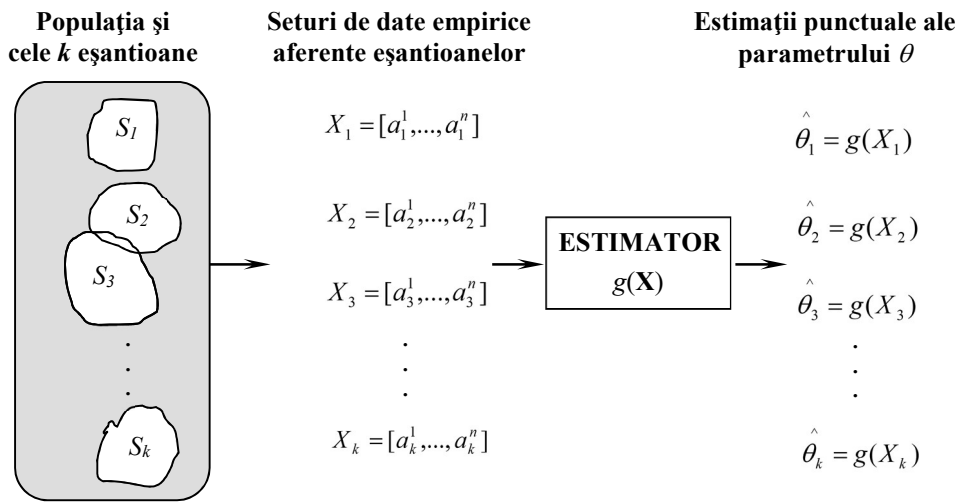
Acești parametri ai modelului statistic cu structură (*fdp*) predefinită nu sunt oarecare ci ei au o semnificație particulară, fiind chiar **parametri ai populației** studiate și având valori tipice (fiind, de exemplu, momente ale distribuției variabilei aleatoare ce caracterizează caracteristica studiată a populației). În cazul nostru particular,  $m_x$  și  $\sigma_x^2$  corespund următorilor doi parametri ai populației studiate (sau, momente de ordin unu, respectiv, doi ale distribuției):

(a) *media* (valoarea sperată a) v.a.  $x$ , care prin definiție este:

$$m_x \stackrel{def}{=} E\{x\} = \int_{-\infty}^{+\infty} af_x(a)da \quad (5.116)$$

(b) și, *varianța* v.a.  $x$ , dată de:

$$\sigma_x^2 \stackrel{def}{=} E\{(x-m_x)^2\} = \int_{-\infty}^{+\infty} (a-m_x)^2 f_x(a)da \quad (5.117)$$



**Figura 5.16.** Procesul de estimare punctuală a parametrului  $\theta$  al unui model statistic parametric

În acest context, funcția densitate de probabilitate din relația (5.115) mai este referită ca *repartiție normală* cu media  $m_x$  și dispersia  $\sigma_x^2$ , valorile celor doi parametri determinând complet repartiția (poziționarea și alura acesteia). Faptul că *v.a.*  $x$  este repartizată normal cu parametrii,  $m_x$  și  $\sigma_x^2$ , se mai notează și:

$$x \sim N(m_x, \sigma_x^2) \quad (5.118)$$

**Estimarea** (*punctuală* sau *de tip interval*) a **parametrilor** populației se face pe câte un eșantion<sup>38</sup> (vezi **Figura 5.16**) – ales din populația investigată astfel încât acesta să fie reprezentativ.

Pentru problematica tratată în această carte de interes particular este estimarea punctuală, despre care vom vorbi în cele ce urmează.

### **Estimarea parametrică punctuală clasică**

Procesul estimării punctuale a parametrilor unei populației sau ai unui proces se implementează cu ajutorul unui **estimator**<sup>39</sup>; acesta folosește la intrare datele măsurate ale eșantionului și furnizează la ieșire *valoarea estimată* (aproximată) a parametrilor necunoscuți ai populației. Exprimat matematic, aceasta s-ar traduce astfel:

- Presupunem că dorim să estimăm parametrul  $\theta$  (scalar sau vectorial) al unei fdp cunoscute,  $f_x(a; \theta)$ , a *v.a.*  $x$  (discuția se poate generaliza la vectori aleatori). Pentru aceasta extragem dintr-o populație  $k$  eșantioane de volum  $n$ , respectiv, eșantioanele  $S_1, \dots, S_k$ . (similar ca în **Figura 5.16**). Fiecare eșantion,  $S_i$ , generează un vector al datelor observabile,  $X_i = [a_i^1, \dots, a_i^n]$ .
- **Estimatorul**  $\hat{\theta}$  al parametrului  $\theta$  al populației este atunci definit ca o funcție  $g(\cdot)$  ce se aplică vectorului aleator de date empirice  $\mathbf{X} = [\mathbf{a}^1, \dots, \mathbf{a}^n]$ <sup>40</sup>; dacă în loc de *v.a.*  $x$  am fi avut vectorul aleator  $x$ , atunci am fi vorbit de o matrice aleatoare de date empirice.

Pentru a înțelege mai bine cele ce urmează, dăm și următoarele definiții pentru o *v.a.*  $h(x)$  a lui  $x$ , pentru media acesteia și pentru noțiunea de statistică:

*Fie  $x$  o v.a. și  $h(x)$  o funcție a v.a.  $x$ . Expresia:*

<sup>38</sup> Submulțime a populației statistice considerate.

<sup>39</sup> Formulă sau procedură de estimare.

<sup>40</sup> Notațiile adiacente folosite în acest subcapitol au fost  $a_i^j$ ,  $\bar{x}$ ,  $\hat{\theta}_i$ ,  $X_i$  = valori particulare ale variabilelor aleatoare/ vectorilor aleatori/ matricelor aleatoare  $\mathbf{a}^j$ ,  $\bar{\mathbf{x}}$ ,  $\hat{\boldsymbol{\theta}}$ , respectiv,  $\mathbf{X}$ .

$$y = h(x) \quad (5.119)$$

este **o nouă v.a. definită astfel**: pentru un  $\zeta$  dat,  $x(\zeta)$  este un număr iar  $h[x(\zeta)]$  – valoarea funcției  $h(\cdot)$  calculată pentru acest număr – este valoarea lui  $y(\zeta) = h[x(\zeta)]$  pentru care mulțimea domeniului este mulțimea  $S$  a realizărilor experimentale. (definiția de mai sus este valabilă și pentru  $x$  vector aleator, cu mențiunea că de această dată  $y$  poate fi v.a., vector aleator sau chiar o matrice aleatoare).

Orice funcție care se aplică vectorului eșantion  $\mathbf{X}$  se mai numește și **statistică**. Astfel, statisticile sunt caracteristici matematice ale eșantioanelor, și ele pot fi utilizate ca valori estimate ale parametrilor (ce sunt, la rândul lor, caracteristici matematice ale populațiilor din care au fost luate eșantioanele).

Din relația (5.116) rezultă că **media v.a.  $y$  este și ea dată de**:

$$E\{y\} = \int_{-\infty}^{+\infty} b f_y(b) db \quad (5.120)$$

Conform unei teoreme de bază (lăsăm ca exercițiu demonstrarea acesteia),  $E\{y\}$  poate fi exprimat nu neapărat în termenii fdp asociate lui  $y$ ,  $f_y(b)$ , ci direct în termenii funcției  $f_x(a)$  a lui  $x$ , ca în relația de mai jos:

$$E\{h(x)\} = \int_{-\infty}^{+\infty} h(x) f_x(a) da \quad (5.121)$$

Din cele de mai sus reiese că și estimatorul punctual  $\hat{\theta} = g(\mathbf{X})$  este tot o v.a. (sau, funcție de situație, un vector aleator/o matrice aleatoare), iar din modul cum a fost definit, el se mai numește și statistică.

- *Estimatul punctual* este rezultatul (valoarea) funcției  $g(\cdot)$  aplicată unui singur eșantion,  $\hat{\theta}_i = g(X_i)$ . În eșantioane diferite, statisticile calculate au valori diferite. Prin urmare, se poate vorbi despre o distribuție a valorilor statisticii  $\hat{\theta}$  în mulțimea eșantioanelor de același volum,  $n$ , numită și *distribuție de sondaj*<sup>41</sup> a statisticii respective.

<sup>41</sup> Operațiunea de formare a unui eșantion (selecție) se numește sondaj.

Din cele prezentate până acum reiese că, atunci când facem inferențe privind o caracteristică studiată din populație operăm, în fapt, cu *trei distribuții asociate acesteia*:

- distribuția populației – este acea distribuție (în general, necunoscută) pe care o are caracteristica studiată (sau *v.a.* asociată ei) în populație și pe care dorim să o determinăm;
- distribuția eșantionului – este distribuția pe care o are *v.a.* în datele empirice (eșantionul) de care dispunem și ea poate fi determinată complet din aceste date;
- distribuția de sondaj – este distribuția pe care o are statistica calculată în mulțimea tuturor eșantioanelor de volum,  $n$ , dat.

Aparent complicat, lucrurile se simplifică însă foarte mult atunci când, datorită unor teoreme de limită centrală, se demonstrează că forma distribuției de sondaj este una cunoscută atunci când volumul eșantionului,  $n$ , crește (vezi **Anexa: Teorema limită centrală**). Mai mult, din considerente teoretice, relația repartiției de sondaj a statisticii cu valoarea tipică (parametrul) din populație este, de asemenea, una bine precizată.

Fără a detalia prea mult din statistica matematică ce stă la baza acestor aspecte teoretice, vom enunța în continuare numai acele rezultate ale teoriei statistice de care ne vom folosi în mod deosebit în această carte, ca parte a statisticii aplicate. Astfel:

- I. Pentru orice parametru al populației,  $\theta$ , dat pot fi implementați unul sau mai mulți estimatori,  $\hat{\theta} = g(\mathbf{X})$ . Obiectivul estimării este acela de a selecta o funcție  $g(X)$  care să minimizeze într-un anumit sens eroarea de estimare  $g(\mathbf{X}) - \theta$ . Atunci când minimizarea se face în sensul erorii pătratice medii, respectiv:

$$e = E\{[g(X) - \theta]^2\} = \int_{\mathfrak{R}} [g(X) - \theta]^2 f_x(X, \theta) dX \quad (5.122)$$

atunci estimatorul  $\hat{\theta} = g(\mathbf{X})$  se numește **cel mai bun estimator**.

- II. În **Subcapitolul** [Error! Not a valid link.](#) s-a făcut mențiunea că *mediile calculate pe un număr de observații se apropie de o valoare constantă pe măsură ce numărul observațiilor crește*. Aceasta reprezintă, în fapt, *interpretarea frecvențială a probabilității*, conform căreia, media

aritmetică  $\bar{x}$  a valorilor observate  $a^i$  ale v.a.  $x$  tinde către integrala din (5.116), ce reprezintă media v.a.  $x$ , atunci când  $n \rightarrow \infty$ :

$$\bar{x} = \frac{a^1 + \dots + a^n}{n} \rightarrow E\{x\} = m_x \quad (5.123)$$

III. În contextul notațiilor de mai sus, v.a.  $\bar{x} = \hat{\theta} = g(\mathbf{X})$ , reprezentând statistica *media aritmetică a eșantionului* a v.a.  $x$ :

$$\bar{x} = \frac{a^1 + \dots + a^n}{n} \quad (5.124)$$

este o v.a. cu media  $m_x$  și varianța  $\sigma_x^2/n$  (mai exact, se acceptă conform unor teoreme de limită centrală (vezi **Anexa: Teorema limită centrală**) că distribuția de sondaj a mediei este caracterizată de o distribuție  $N(m_x, \sigma_x^2/n)$ ); această afirmație este valabilă, în general, pentru valori ale lui  $n > 10$ , atunci când distribuția lui  $x$  este aproape simetrică, și pentru  $n > 30$ , atunci când distribuția lui  $x$  are asimetrie pronunțată sau necunoscută.

Folosindu-ne de toate informațiile prezentate mai sus se poate arăta că statistica *media aritmetică a eșantionului*,  $\bar{x}$ , a v.a.  $x$  este, din punct de vedere teoretic, un:

- *estimator nedeplasat* ( $E\{\hat{\theta}\} = \theta$ ) al parametrului *media*,  $m_x$ , al populației, adică  $E\{\bar{x}\} = m_x$ , și un
- *estimator consistent* (eroarea de estimare  $[\hat{\theta} - \theta]$  tinde către 0 în probabilitate atunci când  $n \rightarrow \infty$ ); mai exact, varianța lui  $\bar{x}$ ,  $\sigma_x^2/n \rightarrow 0$  atunci când  $n \rightarrow \infty$  de unde rezultă că  $\bar{x} \rightarrow m_x$  în sensul mediei pătratice, și, deci, de asemenea, în probabilitate.

**Observația 5.25:** Un estimator am spus că este definit ca o funcție a eșantionului (setului de date empirice) însă, pentru ca el să fie un estimator de încredere este necesar să posede următoarele proprietăți: să fie nedeplasat, să fie consistent și să aibă varianță minimă.

Întrucât decizia de estimator consistent se poate lua doar atunci când  $n \rightarrow \infty$ , ceea ce nu este cazul în practică întrucât  $n$  poate lua cel mult valori finite, rezultă că cel mult putem căuta și, în cel mai bun caz, selecta empiric, cel mai bun estimator.

În cele ce urmează, alegerea empirică a estimatorului o vom face ținând cont de două lucruri:

1. statistica *media aritmetică a eșantionului* a lui  $x$ , din relația (5.124),



$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n a^i \quad (5.125)$$

este, așa după cum am văzut din punct de vedere teoretic, un estimator consistent al parametrului  $\theta = m_x = E\{x\}$ , ce reprezintă media populației pentru caracteristica studiată a populației;

2. din definiția dată în relația (5.119) avem că pentru o v.a.  $x$ , o funcție  $h(x) = y$  este tot o v.a.

Combinând cele două informații ajungem la **construirea empirică a următorului estimator** pentru parametrul *media*  $\theta = E\{h(x)\}$  ( $= m_y = E\{y\}$ ) a funcției  $h(x)$  a lui  $x$ , respectiv:

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n h(a^i) \quad (5.126)$$

Acest estimator reprezintă media aritmetică a eșantionului pentru v.a.  $y$ , deci pentru funcția  $h(x)$ , și el este:

- din *punct* de vedere teoretic, un estimator consistent iar
- din punct de vedere practic, de cele mai multe ori, cel mai bun estimator (asta cel puțin pentru  $n$  de valori mari).

Funcție de forma pe care o îmbracă funcția  $h(x)$  identificăm, de exemplu, următoarele cazuri particulare:

a) $h(x)=x$ :	$\theta = E\{x\} = m_x$ și $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n a^i$	→ media aritmetică a eșantionului este estimator pentru media populației;
b) $h(x)=(x-m_x)^2$ :	$\theta = E\{(x-m_x)^2\} = \sigma_x^2$ și $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n (a^i - m_x)^2$	→ varianța eșantionului este estimator pentru varianța populației;
c) <i>etc.</i>		

unde, **parametrul  $\theta$  al populației are**, așa cu cum se poate remarca, **valori tipice**, respectiv, **valori asociate diverselor momente statistice ale v.a.  $x$** .

**Observația 5.26:** Așa după cum am prezentat și în **Subcapitolul 5.2, elemente fundamentale de statistică** – spre deosebire de *variabilele statistice* ce desemnează o anumită caracteristică a populației –, prin parametru al

aceleiași populației am spus că se înțelege doar o proprietate numerică a acesteia. În acest moment, cu informațiile prezentate până acum, suntem în măsură să dăm următorul exemplu simplu care să ne ajute să înțelegem (în caz că nu am reușit să o facem deja) distincția dintre cele două noțiuni. Astfel, dacă luăm în considerare parametrul *media* (valoarea sperată) a unei *v.a.*  $x$  (fie aceasta de tip continuu) ce descrie o caracteristică a populației, avem că ea este teoretic dată de relația (5.116):

$$m_x = E\{x\} \stackrel{\text{def}}{=} \int_{-\infty}^{+\infty} af_x(a)da \quad (5.127)$$

iar empiric, ea este estimată prin media aritmetică a eșantionului (relația (5.123)):

$$\bar{x} = \frac{a^1 + \dots + a^n}{n} \quad (5.128)$$

care, la limită, pentru  $n \rightarrow \infty$ , tinde spre valoarea reală a parametrului  $m_x$ . Cu alte cuvinte, atunci când ea există, *media* lui  $x$  este, aproape sigur, limita mediei eșantionului obținută pentru o mărime a eșantionului ce crește la infinit.

În acest context, termenul echivalent de “*valoare sperată*” dat parametrului *media* este unul care ne poate induce în eroare, el putând fi, în mod eronat, asociat cu „*valoarea cea mai probabilă*” a respectivei variabile statistice. La fel ca și pentru media eșantionului, și pentru media populației putem obține însă o valoare „*imposibilă*” în sensul că ea nu reprezintă o valoare tipică dintre valorile posibile pe care le poate lua efectiv *v.a.*  $x$ .

Astfel, dacă, de exemplu, considerăm eșantionul reprezentat de totalitatea pacientelor dintr-un spital de obstetrică-ginecologie, pentru care observăm statistic variabila statistică  $x$  ce desemnează numărul de nașteri realizate de către fiecare pacientă (unitate statistică), atunci un estimat al mediei acestei variabile statistice, dat de media aritmetică a eșantionului, s-ar putea să ne conducă la un rezultat de tipul **3.7** nașteri/pacientă; acest rezultat statistic, care nu poate reprezenta nicidecum o valoare posibilă în mod real, ne dă doar o informație orientativă și anume aceea că, în medie, fiecare pacientă a avut circa 3.7 nașteri până la data realizării studiului.

**Observația 5.27:** O abordare similară a problemei de estimare parametrică punctuală o întâlnim și în cazul în care avem de-a face cu **statistici**

**comune** pentru două variabile aleatoare,  $x$  și  $y$ , ce caracterizează două trăsături distincte ale populației. Aceasta implică următoarele:

- date două variabile aleatoare,  $x$  și  $y$ , și o funcție  $h(a,b)$ , formăm v.a.  $z = h(x,y)$ ;
- parametrul media  $\theta = E\{h(x,y)\} = E\{z\}$  a funcției  $h(x,y)$  a variabilelor aleatoare  $x$  și  $y$  poate fi exprimat și el direct în termenii funcției  $h(x,y)$  și ai funcției densitate comună,  $f_{xy}(a,b)$ , astfel:

$$E\{h(x,y)\} = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} h(a,b) f_{xy}(a,b) da db \quad (5.129)$$

- se poate, de asemenea, demonstra că un *estimator empiric*,  $\hat{\theta}$ , al parametrului  $\theta$  de mai sus este statistica *media aritmetică a eșantionului* pentru v.a.  $z$ , deci pentru funcția  $h(x,y)$ , a lui  $x$  și  $y$ , respectiv:

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n h(a^i, b^i) \quad (5.130)$$

Funcție de forma particulară pe care o îmbracă funcția  $h(x,y)$  identificăm, de exemplu, următoarele cazuri, deosebit de utile în aplicațiile practice, pentru parametrul  $\theta$  și estimatorul său,  $\hat{\theta}$ :

a) $h(x,y) = xy$ :	$\theta = E\{xy\}$ $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n a^i b^i$
b) $h(x,y) = (x-m_x)(y-m_y)$ :	$\theta = E\{(x-m_x)(y-m_y)\} = cov(x,y)$ $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n (a^i - m_x)(b^i - m_y)$
c) <i>etc.</i>	

- în cazul particular când variabilele aleatoare  $x$  și  $y$  au o distribuție comună normală funcția densitate comună,  $f_{xy}(a,b)$  este determinată în mod unic de următoarele momente de ordinul întâi și doi:  $m_x$ ,  $m_y$ ,  $\sigma_x$ ,  $\sigma_y$  și  $cov(x,y)$  (vezi și **Subcapitolul 5.6**).

Până în acest punct al prezentării noastre parametrii populației au fost tratați ca fiind niște *constante necunoscute* – aceasta reprezintă, de fapt, ipoteza de lucru de bază a *estimării clasice*. În estimarea parametrică mai există însă și o abordare în care parametrii populației sunt ei însăși văzuți nu ca niște cantități necunoscute și *fixe* ci ca niște *variabile aleatoare*, având

propriile lor distribuții de probabilitate. Această nouă abordare poartă numele de *estimare bayesiană* a parametrilor populației.

### *Estimarea parametrică punctuală Bayes-iană*

În estimarea Bayes-iană a parametrilor unei populații se pleacă de la următoarele *două premise*:

- (1) parametrul  $\theta$  al populației nu este în totalitate unul necunoscut (altfel spus, dispunem de anumite informații *a priori* fie dintr-o cunoaștere parțială a procesului studiat, fie din alte realizări anterioare ale aceluiași experiment etc.), și
- (2) întrucât valoarea parametrului este necunoscută atunci are sens să specificăm o distribuție de probabilitate care să descrie valorile posibile pentru acest parametru, la fel ca și probabilitățile asociate acestor valori.

Informațiile *a priori* despre parametrul  $\theta$  sunt valorificate în momentul în care construim și precizăm funcția de densitate *a priori*,  $f_{\theta}(\theta)$ , a vectorului aleator  $\theta$ . În această abordare, parametrul necunoscut al populației este doar o valoare a vectorului aleator  $\theta$  iar funcția de distribuție a vectorului aleator  $x$ , ce descrie caracteristica populației, este interpretată ca fiind funcția de distribuție condiționată,  $F_x(a|\theta)$ , a lui  $x$  pentru  $\theta = \theta$ .

În aceste ipoteze de lucru, *metoda bayesiană de estimare statistică* apare ca o alternativă la *metoda clasică de estimare statistică* de care se diferențiază, însă, în principal, prin următoarele două elemente cheie:

- i) parametrul populației este el însuși văzut ca un vector aleator,  $\theta$ , valoarea sa necunoscută, fiind o valoare a acestui vector aleator;
- ii) informația furnizată de setul de date empirice, de care dispunem la un moment dat, este combinată, în plus, cu informația a priori pe care o avem despre parametrul  $\theta$  al populației. Reamintim că, în estimarea parametrică clasică, informația folosită în ghidarea procesului de inferență statistică era doar cea furnizată de setul de date empirice.

O concluzie imediată ce poate fi trasă din cele de mai sus este aceea că în abordarea Bayes-iană, **problema estimării parametrului necunoscut al populației revine la a estima valoarea  $\theta$  a lui  $\theta$  în termenii:**

- eșantionului  $X$  (constând în valorile observate  $a^i$  ale lui  $x$ ) și, respectiv,
- funcției de densitate de probabilitate a lui  $\theta$ .

**Observația 5.28:** Atenție! Prin transformarea problemei dintr-o estimare a parametrului necunoscut  $\theta$  (întâlnită în estimarea clasică) într-o problemă de estimare a valorii  $\theta$  a vectorului aleator  $\boldsymbol{\theta}$  (în cazul abordării bayes-iene), **estimarea** devine în această ultimă abordare, practic, **o predicție** (vezi diferența dintre acestea în **Anexa: Predicție versus estimare**).

Ținând cont de faptul că problema de estimare s-a transformat într-o problemă de predicție (vezi **Observația 5.28**, respectiv, relația (A.71)), atunci putem afirma că soluția pe care o căutăm – fie aceasta constanta  $\hat{\theta}$  :

- **era**, înainte de a obține orice evidență statistică (date observabile), **dată de**,

$$\hat{\theta} = E\{\boldsymbol{\theta}\} = \int_{-\infty}^{+\infty} \boldsymbol{\theta} f_{\boldsymbol{\theta}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (5.131)$$

- iar după obținerea de noi probe (respectiv, obținerea unui set de date empirice pentru vectorul aleator  $x$ ,  $X=[a^1, \dots, a^n]$ ), ea **este dată de** (vezi relația (5.187)),

$$\hat{\boldsymbol{\theta}} = E\{\boldsymbol{\theta} | X\} = \int_{-\infty}^{+\infty} \boldsymbol{\theta} \cdot f_{\boldsymbol{\theta}|X}(\boldsymbol{\theta} | X) d\boldsymbol{\theta} \quad (5.132)$$

În formula de mai sus  $f_{\boldsymbol{\theta}|X}(\boldsymbol{\theta} | X)$  este tot funcția densitate de probabilitate a lui  $\boldsymbol{\theta}$  însă revizuită în sensul că funcția de densitate *a priorică*,  $f_{\boldsymbol{\theta}}(\boldsymbol{\theta})$  – cunoscută în totalitate și construită în baza tuturor informațiilor de care dispuneam înainte de a obține noi rezultate experimentale – este modificată corespunzător, în contextul și conform noilor dovezi statistice obținute experimental.

Derivarea noii funcții de densitate pentru vectorul aleator  $\boldsymbol{\theta}$  are ca regulă de calcul regula lui Bayes (vezi relația (5.93)),

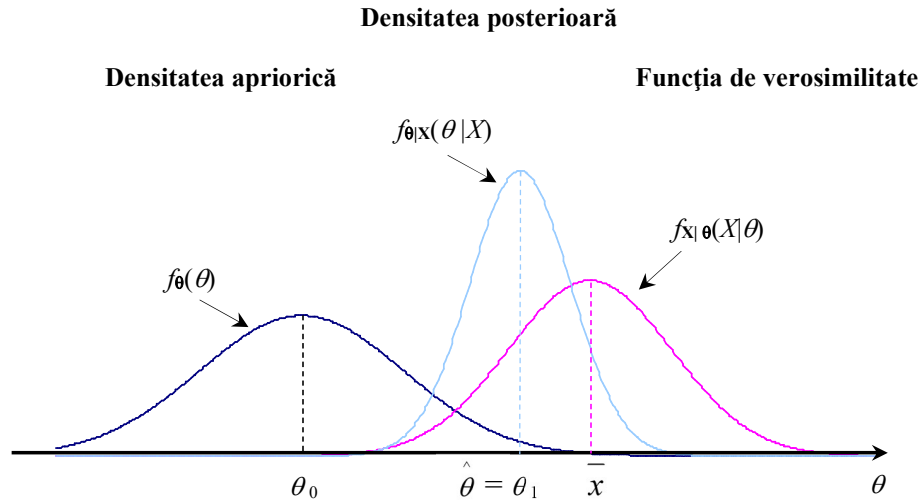
$$f_x(a | B) = \frac{P(B | x = a) f_x(a)}{P(B)} = \frac{P(B | x = a) f_x(a)}{\int_{-\infty}^{+\infty} P(B | x = a) f_x(a) da} \quad (5.133)$$

particularizată pentru  $x = \boldsymbol{\theta}$  și, respectiv, pentru evenimentul  $B = \mathbf{X}$ , ce constă în obținerea unui eșantion de dimensiune  $n$  pentru vectorul aleator  $x$ . Astfel:

$$f_{\boldsymbol{\theta}|X}(\boldsymbol{\theta} | X) = \frac{f_{X|\boldsymbol{\theta}}(X | \boldsymbol{\theta}) f_{\boldsymbol{\theta}}(\boldsymbol{\theta})}{f_X(X)} \quad (5.134)$$

În relația de mai sus notațiile folosite sunt aceleași cu cele folosite în estimarea parametrică clasică (vezi **Subcapitolul 5.5.1.**), respectiv:

- $\theta$  este vectorul aleator ce desemnează parametrul populației,
- $\theta$  este o realizare particulară a acestui vector aleator,
- $\mathbf{X}$  desemnează vectorul aleator sau matricea aleatoare (după cum  $x$  este un scalar, respectiv, un vector) a eșantionului vectorului aleator  $x$  (vezi **Figura 5.17**) iar
- $X$  reprezintă o realizare particulară a lui  $\mathbf{X}$ , respectiv, setul de date empirice pe baza căruia se realizează inferența statistică.



**Figura 5.17.** Exemplu de estimare bayesiană punctuală parametrică în ipoteza unor distribuții normale ale lui  $\theta$  și, respectiv,  $x$ .

Această formulă de calcul, ilustrată grafic în **Figura 5.17**, ne arată cum, pornind de la informația disponibilă – specificată prin funcția de densitate a priori,  $f_{\theta}(\theta)$ , pe care o presupunem cunoscută –, în lumina unei noi evidențe statistice – respectiv, datele observabile ale eșantionului, a căror informație furnizată este reflectată în funcția de verosimilitate,  $f_{X|\theta}(X|\theta)$  – derivăm o funcție de densitate de probabilitate actualizată pentru vectorul aleator  $\theta$ ,  $f_{\theta|X}(\theta|X)$ , funcție denumită, în mod sugestiv, funcție densitate de probabilitate posterioară sau funcție de densitate de probabilitate condiționată a lui  $\theta$  dată de evidența  $X$ .

**Observația 5.29.:** Facem observația aici că pentru funcția  $f_{X|\theta}(X|\theta)$  – ce reprezintă densitatea condiționată comună a  $n$  vectori aleatori,  $\mathbf{a}^i$ , atunci

când  $\hat{\theta} = \theta$  –, este adevărată următoarea relație de calcul pentru cazul când acești vectori aleatori sunt independenți:

$$f_{X|\theta}(X|\theta) = f_{a^1|\theta}(a^1|\theta) \cdots f_{a^n|\theta}(a^n|\theta) \quad (5.135)$$

Aici, funcția generică  $f_{a|\theta}(a|\theta)$  reprezintă funcția densitate condiționată a vectorului aleator  $\mathbf{a}$  atunci când  $\theta = \theta$ .

În situația particulară când funcția  $f_{X|\theta}(X|\theta)$  este considerată ca o funcție de  $\theta$ , ea poartă numele de funcție de verosimilitate; acesta este și cazul funcției  $f_{X|\theta}(X|\theta)$  din relația (5.198), unde pentru o valoare particulară a lui  $\mathbf{X}$  (un eșantion  $X_k$  oarecare, vezi **Figura 5.16**) avem:

$$f_{X|\theta}(X_k|\theta) = f_{a^1|\theta}(a^1|\theta) \cdots f_{a^n|\theta}(a^n|\theta) \quad (5.136)$$

**Problema 5.22:** Să se estimeze punctual, prin metoda *bayes-iană*, parametrul scalar *media*,  $\theta = m_x$ , al unei populații a cărei caracteristică studiată este descrisă de vectorul aleator unidimensional,  $x$ , despre care se știe că este distribuit conform unei legi de distribuție normale,

$$x \sim N(\theta, \sigma_x^2) \quad (5.137)$$

și pentru care se cunoaște doar varianța,  $\sigma_x^2$ . De asemenea, mai dispunem:

- de o serie de informații *a priori* despre parametrul  $\theta$  care ne fac să credem că valorile posibile pentru acesta respectă, de asemenea, o distribuție normală, de medie și varianță, cunoscute, respectiv,

$$\theta \sim N(\theta_0, \sigma_0^2) \quad (5.138)$$

- și de un eșantion, de dimensiune  $n$ , de date empirice.

Să se comenteze rezultatul obținut.

*Notă* Se arată că estimatul  $\hat{\theta} = \theta_1$  pentru parametrul necunoscut al populației este unul dat de formula:

$$\theta_1 = \frac{n\sigma_1^2}{\sigma_x^2} \bar{x} + \frac{\sigma_1^2}{\sigma_0^2} \theta_0 \quad (5.139)$$

unde  $\bar{x}$  reprezintă statistica media aritmetică a eșantionului pentru v.a.  $x$  iar,

$$\sigma_1^2 = \frac{\sigma_x^2 \sigma_0^2}{\sigma_x^2 + n \sigma_0^2} \quad (5.140)$$

**Observația 5.30:** În problema de mai sus se poate remarca faptul că, atât distribuția parametrului  $\theta$ , cât și distribuția caracteristicii populației (reprezentată prin *v.a.*  $x$ ), sunt distribuții normale iar distribuția posterioară obținută pentru  $\theta$  în urma aplicării regulei lui *Bayes* este, de asemenea, o distribuție normală. Rezultatul nu este, așa după cum s-ar putea crede, unul total întâmplător. Mai exact, dacă în general în estimarea parametrică bayesiană familia de distribuții căreia îi aparține funcția de verosimilitate,  $f_{X|\theta}(X|\theta)$ , se consideră a fi una bine-determinată (fixată), în ceea ce privește distribuția *a priori* a parametrului populației,  $f_{\theta}(\theta)$ , avem doar presupunerea crucială că aceasta este o distribuție cunoscută, rămânând ca alegerea uneia sau a alteia dintre familiile de distribuții teoretice existente să se facă, în parte pe baza informațiilor *a priori* disponibile, iar în parte, pe baza intuiției noastre. Este însă evident faptul că alegeri diferite pentru densitatea *a priori*  $f_{\theta}(\theta)$  poate face ca produsul  $f_{X|\theta}(X|\theta) \times f_{\theta}(\theta)$  să ia o formă algebrică arbitrară sau alta (de cele mai multe ori, destul de greu de descris), formă care îngreuiază foarte mult, dacă nu chiar, face imposibil calculul distribuției posterioare,  $f_{\theta|X}(\theta|X)$ .

Soluția în astfel de situații o reprezintă folosirea așa – numitelor **densități *a priori* conjugate**.

O clasă de distribuții de probabilitate *a priori*,  $f_{\theta}(\theta)$ , se numește **conjugata** unei clase de funcții de verosimilitate,  $f_{X|\theta}(X|\theta)$ , doar atunci când distribuțiile posterioare rezultante,  $f_{\theta|X}(\theta|X)$ , fac parte din aceeași familie ca și  $f_{\theta}(\theta)$ . Practic, o probabilitate *a priori* conjugată este o conveniență algebrică care simplifică foarte mult scrierea matematică, oferind, în acest mod o formulă simplă de calcul ce leagă valorile hiper-parametrilor<sup>42</sup> densității posterioare de valorile hiper-parametrilor densității *a priori* [Gelman, 2003], [Fink, 1995]. În acest fel, și calculul densității posterioare devine unul foarte ușor, densitatea posterioară fiind, în final, o densitate din aceeași familie de distribuție ca și densitatea *a priori* însă cu hiperparametri diferiți ce reflectă efectul, suplimentar, al informației furnizate de datele empirice.

<sup>42</sup> Un **hyper-parametru** reprezintă un parametru al distribuției *a priori* iar termenul este utilizat pentru a distinge acest parametru de parametrii modelului propus pentru distribuția ce descrie caracteristica populației.



Revenind la **Problema 5.22** putem acum bănuî că presupunerea conform căreia distribuția *a priori* pentru  $\theta_0$  este o distribuție normală a fost una de natură să simplifice scrierea matematică. Uitându-ne în tabelul cu funcții densitate *a priori* conjugate, prezentat în **Anexa: Distribuții de verosimilitate**, observăm că, într-adevăr, familia de distribuții *gauss*-iene este conjugată cu ea însăși (auto-conjugată). Mai mult, corespunzător celor prezentate mai sus și ținând cont de soluția dată problemei identificăm:

Funcția de verosimilitate	Parametrii necunoscuți ai modelului	Distribuția <i>a priori</i> conjugată	Hyper-parametrii <i>a priori</i>	Hyper-parametrii posteriori
Normală cu varianță, $\sigma_x^2$ , cunoscută	$\theta = m_x$ (media)	Normală	$\theta_0, \sigma_0^2$	$\theta_1 = \left( \frac{n\sigma_0^2}{\sigma_x^2 + n\sigma_0^2} \right) \bar{x} + \frac{\sigma_x^2}{\sigma_x^2 + n\sigma_0^2} \theta_0$ $\sigma_1^2 = \frac{\sigma_x^2 \sigma_0^2}{\sigma_x^2 + n\sigma_0^2}$

### Probabilitatea *a priori*

Deși subiectivă în esență (a se vedea și **Anexa: Interpretarea noțiunii de probabilitate**), probabilitatea *Bayes*-iană este interpretată uneori ca fiind **obiectivă**. Cele două puncte de vedere diferite – varianta **subiectivă** și, respectiv, varianta **obiectivă** ale probabilității bayesiene – diferă, în principal, prin modul cum este **interpretată și construită probabilitatea *a priori*** din relația lui Bayes. Astfel:

- conform punctului de vedere **subiectiv**, probabilitatea este interpretată ca fiind “*gradul de încredere*” pe care un individ îl are în adevărul unui enunț/eveniment; de aici rezultă și caracterul personal sau subiectiv; cu alte cuvinte, nivelul de cunoaștere corespunde unei “convingeri personale”;
- conform punctului de vedere **obiectiv**, nivelul *a priori* al cunoașterii definește – pentru problemele bine puse – o distribuție de probabilitate *a priori* unică. În general, problemele de inferență statistică, așa cum se întâlnesc ele în investigarea științifică, sunt rezolvate, în mod uzual, plecând de la presupunerea că un anume model statistic este un descriptor adecvat al mecanismului probabilistic ce a generat datele; însăși alegerea

acelui model implică, însă, în continuare, un anumit grad de subiectivitate. Cu toate acestea, a devenit o practică standard ca orice analiză statistică ce depinde doar de modelul propus și de datele observate să fie descrisă ca fiind “obiectivă”. În acest sens (și doar în acest sens) inferența bayesiană poate fi privită ca fiind “obiectivă”. În construirea obiectivă a distribuțiilor *a priori*, au fost propuse utilizarea mai multor principii dintre care amintim doar: *entropia maximă* (Edwin T. Jayne), *analiza de referință* (Jose M. Bernardo) și *analiza frecvențială* (în această ultimă situație distribuția de frecvență are parametri ce reflectă o analiză și o sinteză a unor date existente).

În ceea ce privește **densitățile *a priori* de referință**, acestea nu sunt, așa cum s-ar putea crede, descriptori ai gradelor subiective de încredere; ele sunt propuse ca funcții *a priori* formale, unanim acceptate spre a fi utilizate ca standarde în comunicările științifice.

Ceea ce trebuie să reținem în final, de aici, este faptul că **multe dintre metodele moderne de învățare mașină și de clasificare se bazează, în general, pe principiile bayesiene obiective.**

### **Estimarea parametrică a funcției de regresie**

Cele două metode de estimare parametrică discutate până acum reprezintă, așa după cum știm *metode de estimare a funcțiilor densitate de probabilitate* aferente diferitelor variabile aleatoare (mai general, vectori aleatori). Reluând puțin, vom spune că, în cadrul acestor metode se pleacă de la un eșantion,  $X$ , de date observate ale variabilei aleatoare  $x$ , se presupune o anumită formă pentru densitatea de probabilitate a lui  $x$ ,  $f_x(a)$  (suplimentar, în cazul estimării Bayes-iene, se precizează și densitatea *a priorică* a parametrului  $\theta$  al modelului,  $f_\theta(\theta)$ ) și, pe baza acestora, se estimează parametrul  $\theta$ , ce implicit determină complet estimatul funcției densitate,  $f_x(a)$ , a lui  $x$ .

Un alt tip de aplicații statistice cu care ne întâlnim destul de frecvent în practică și care implică, de asemenea, o estimare parametrică – de această dată, însă, a unei funcții numită de regresie – este și următorul, în care: se cunoaște un eșantion,  $X$ , de perechi de date ( $a^i, b^i$ ), ale unei perechi de variabile aleatoare ( $x, y$ ), între care știm că există o *relație funcțională*, se cunoaște forma acestei relații funcționale și, pe baza acestora, se estimează vectorul parametru,  $\theta$ , ce descrie complet relația funcțională de mai sus. În acest caz particular, vorbim de o estimare parametrică a unei funcții de regresie. Funcția de regresie are, după cum vom vedea, o utilitate foarte largă, ea putând fi utilizată în aplicații diverse, începând cu aplicațiile de predicție,

inferență statistică, testare de ipoteze și terminând, nu în cele din urmă, cu aplicațiile de modelare a relațiilor cauzale.

**Estimarea unei funcții de regresie** constă în găsirea celui model care aproximează cel mai bine relația funcțională care există între o variabilă aleatoare dependentă,  $y$ , (numită și *variabilă de răspuns*) și, respectiv, una sau mai multe variabile aleatoare independente (numite și predictorii). Formal, ecuația de regresie pentru cazul unei singure *v.a.* independente,  $x$  (generalizarea la mai mulți predictorii este una imediată), se scrie astfel:

$$b = \phi(a; \theta) + e \quad (5.141)$$

În această relație *v.a.* dependentă  $y$  este modelată ca o funcție de predictorul  $x$ , de vectorul parametru corespunzător  $\theta$  (ce reprezintă un vector de componente „constante”) la care se adaugă un termen de eroare ce reprezintă acea variație din  $y$  ce nu a putut fi explicată (modelată). Termenul de eroare este aici tratat ca o variabilă aleatoare. Reamintim că variabilele  $a$  și  $b$  desemnează variabilele generice aferente variabilelor aleatoare  $x$  și, respectiv,  $y$ .

Înainte de a merge mai departe cu prezentarea acestui tip de estimare, vom introduce în cele ce urmează câteva noțiuni noi. Scopul este acela de a ușura înțelegerea modului cum vom obține un estimator pentru funcția de mai sus – funcție numită și *funcție* sau *linie de regresie*.

#### **Media condiționată a unei variabile aleatoare. Funcția mediei condiționate**

Prin **definiție**, **media condiționată** a unei variabile aleatoare  $y$ , presupunând evenimentul  $A$ , este dată de integrala:

$$E\{y | A\} \stackrel{def}{=} \int_{-\infty}^{+\infty} b f_{y|A}(b | A) db \quad (5.142)$$

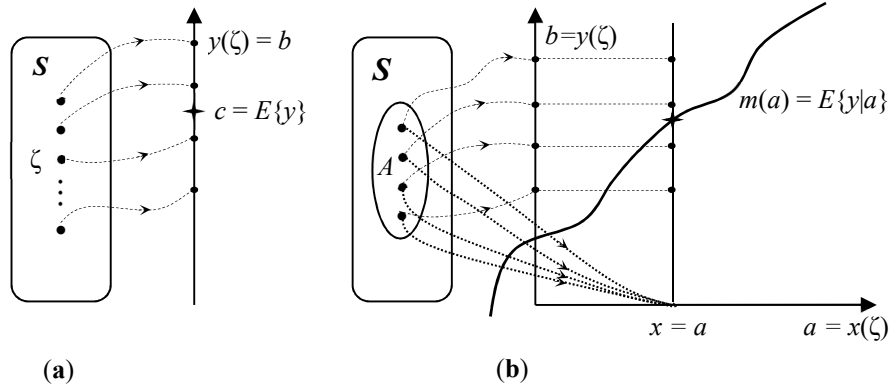
În contextul problemei de estimare de mai sus vom considera pentru evenimentul  $A$  din relația (5.142), următorul eveniment  $\{x = a\}$ . Ținând cont și de relația (5.78) obținem:

$$E\{y | x = a\} = \int_{-\infty}^{+\infty} b f_{y|x}(b | a) db = \int_{-\infty}^{+\infty} b \frac{f_{xy}(a, b)}{f_x(a)} db = m(a) \quad (5.143)$$

unde  $f_{y|x}(b|a)$  este funcția densitate de probabilitate condiționată a lui  $y$  dată de  $x = a$ ,  $f_y(b)$  este *fdp* marginală a lui  $y$  iar  $m(a)$  este o notație pe care o folosim pentru a desemna funcția rezultantă ce depinde doar de variabila  $a$ .

Dacă în relația de mai sus baleiem toată mulțimea valorilor posibile,  $a$ , pentru *v.a.*  $x$  atunci obținem o funcție de la  $a$  la  $b = m(a)$ , numită **funcția mediei condiționate** sau **funcția de regresie**; această funcție, desemnată adesea prin  $E\{y|x\}$ , ne spune cum sunt legate între ele „în medie” variabilele aleatoare  $x$  și  $y$  (vezi **Figura 5.18**). Pentru a înțelege reprezentarea din **Figura 5.18.(a)** vezi problema de predicție tratată în **Anexa: Predicție versus estimare**, și rezumată prin formula (A.71). O înțelegere completă (și nu doar intuitivă) a reprezentării din **Figura 5.18.(b)** o vom obține abia după

prezentarea teoriei *estimării neliniare*, în sensul erorii medii pătratice minime (*mean square - MS*), pe care o redăm în acest subcapitol.



**Figura 5.18.** Reprezentarea grafică a estimării variabilei aleatoare  $y$ : **(a)** printr-o constantă  $c$  și, respectiv, **(b)** printr-o funcție  $m(x)$  a unei v.a.  $x$ .

Funcția mediei condiționate poate fi o funcție neliniară chiar și pentru densități comune,  $f_{xy}(a,b)$ , și densități marginale aparent simple la prima vedere.

**Problema 5.23:** Fie următoarea funcție densitate de probabilitate comună pentru variabilele aleatoare  $x$  și  $y$ :

$$f_{xy}(a,b) = \begin{cases} a + b & , \text{pentru } a \in [0,2] \text{ și } b \in [0,2] \\ 0 & , \text{în rest.} \end{cases} \quad (5.144)$$

Să se deducă funcția mediei condiționate pentru v.a.  $y$ , funcție de v.a.  $x$ .

**Rezolvare:**

Din definiția lui  $f_{xy}(a,b)$  obținem, conform relației (5.43), funcția densitate de probabilitate marginală a lui  $x$  astfel:

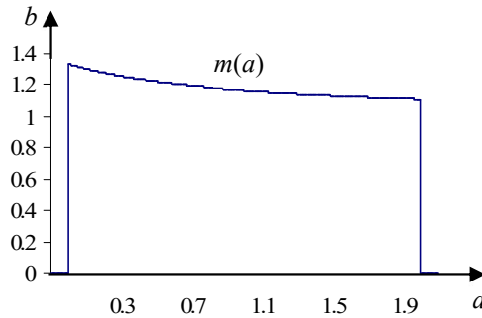
$$f_x(a) = \begin{cases} \int_0^2 f_{xy}(a,b) db = \left( ab + \frac{1}{2} b^2 \right) \Big|_0^2 = 2a + 2 & , \text{pentru } a \in [0,2] \\ 0 & , \text{în rest.} \end{cases} \quad (5.145)$$

Pe baza relației (5.207) deducem mai departe funcția mediei condiționate a lui  $y$  funcție de  $x$ :

$$\begin{aligned} E\{y|x=a\} &= \int_0^2 b \frac{f_{xy}(ab)}{f_x(a)} db = \int_0^2 b \frac{a+b}{2a+2} db = \int_0^2 \frac{a}{2a+2} b db + \int_0^2 \frac{1}{2a+2} b^2 db \\ &= \frac{a}{2a+2} \cdot \frac{b^2}{2} \Big|_0^2 + \frac{1}{2a+2} \cdot \frac{b^3}{3} \Big|_0^2 = \frac{a}{a+1} + \frac{4/3}{a+1} = \\ &= \frac{3a+4}{3a+3} \end{aligned} \quad (5.146)$$

Așa cum arată și rezultatul de mai sus, funcția obținută în acest caz este, evident, o funcție neliniară (vezi figura **Figura 5.19**); acest lucru vine și confirmă afirmația că, chiar și pentru forme simple ale densităților comune și marginale, funcția  $m(a)$  poate fi o funcție neliniară.

Revenind acum la problema estimării relației funcționale,  $\phi(\cdot)$ , ce leagă între ele cele două variabile aleatoare,  $x$  și  $y$ , prezentăm în continuare estimatorul pe care-l obținem în estimarea neliniară în sensul MS a acestei funcții.



**Figura 5.19.** Funcția mediei condiționate,  $m(a)$ .

### *Estimarea neliniară în sensul erorii medii pătratice minime*

Această problemă de estimare revine la a găsi acea funcție  $\phi(\cdot)$  pentru care eroarea medie pătratică:

$$e = E\{[y - \phi(x)]^2\} = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} [b - \phi(a)]^2 f_{xy}(a, b) da db \quad (5.147)$$

este minimă.

Pentru a obține estimatorul corespunzător pornim de la relația de mai sus în care rescriem densitatea comună (vezi și relația (5.69)) astfel:

$$\begin{aligned} e &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} [b - \phi(a)]^2 f_{y|x}(b|a) f_x(a) da db = \\ &= \int_{-\infty}^{+\infty} f_x(a) \int_{-\infty}^{+\infty} [b - \phi(a)]^2 f_{y|x}(b|a) db da \end{aligned} \quad (5.148)$$

În relația (5.212) toți termenii care se integrează sunt pozitivi ceea ce înseamnă că, eroarea  $e$  astfel calculată va fi minimă doar dacă integrala,

$$\int_{-\infty}^{+\infty} [b - \phi(a)]^2 f_{y|x}(b|a) db \quad (5.149)$$

este minimă pentru fiecare  $a$ . Soluția care se obține în acest caz este (a se vedea **Anexa: Predicție versus estimare**, relația (A.68), în care  $c$  este

înlocuit cu  $\phi(a)$  iar  $f_x(a)$  este înlocuit cu  $f_{y|x}(b,a)$ :

$$\phi(a) = \int_{-\infty}^{+\infty} b f_{y|x}(b|a) db = E\{y|a\} = m(a) \quad (5.150)$$

Din relația de mai sus reiese că, estimatul obținut pentru funcția ce leagă între ele cele două variabile aleatoare este chiar *funcția mediei condiționate a variabilei dependente y dată de variabila independentă x*. Cu alte cuvinte, rescriind relația **Error! Not a valid link.** în lumina noilor rezultate obținute, deducem:

$$b = m(a; \theta) + e \quad (5.151)$$

Numele acestei metode de estimare, respectiv, acela de *estimare parametrică neliniară* a regresiei se datorează faptului că forma funcției de regresie este una ce se presupune a fi cunoscută, ea este dată de o combinație neliniară a parametrilor modelului și, în plus, parametrii modelului au o semnificație particulară (uneori chiar fizică) pentru procesul studiat. Metoda estimează acești parametri în sensul găsirii acelor valori pentru care se obține cea mai bună aproximare a datelor observate în condițiile îndeplinirii unui criteriu de eroare anterior stabilit. Estimarea parametrilor se face printr-o *metodă a aproximărilor successive* și, din acest motiv, regresia neliniară este cunoscută ca fiind una intens computațională.

**Observația 5.31:** Funcția de regresie din **Problema 5.23**,

$$m(a) = 1 + \frac{1}{3a+3} = \alpha_0 + \frac{\alpha_1}{\alpha_2 a + \alpha_3} \quad (5.152)$$

este o funcție neliniară întrucât ea nu poate fi exprimată ca o combinație liniară a parametrilor  $\alpha_i$ . Exemple de funcții neliniare mai uzual folosite sunt: funcțiile gaussiene, funcțiile putere, funcțiile exponențiale, funcțiile trigonometrice etc.

O alternativă mai simplă dar, în general, mai puțin eficientă decât metoda parametrică neliniară de estimare a regresiei este și metoda estimării parametrice liniare, în sens MS, a acesteia.

### ***Estimarea liniară în sensul erorii medii pătratice minime***

În cazul liniar univariat, problema estimării parametrice a funcției de regresie dintre variabilele aleatoare  $x$  și  $y$  se reduce la estimarea, în sensul erorii medii pătratice, a *v.a.*  $y$  în termenii unei funcții liniare a lui  $x$ , funcție de tipul  $Ax+B$ . Așa după cum vom vedea mai jos, parametrii cei mai buni de aproximare au, în estimarea liniară a regresiei, expresii bine determinate; în

cazul nostru vom arăta că parametrii estimați se calculează pe baza momentelor de ordin unu și doi ale distribuțiilor lui  $x$  și, respectiv,  $y$ .

Pentru aceasta, scriem, mai întâi formula pentru eroarea medie pătratică, ce trebuie minimizată:

$$e = E\{[y - (Ax + B)]^2\} \quad (5.153)$$

Pentru un  $A$  fixat eroarea  $e$  de mai sus este minimă numai dacă:

$$\begin{aligned} \frac{de}{dB} = 0 &\Leftrightarrow \frac{E\{y^2 - 2(Ax + B)y + [(Ax)^2 + 2AxB + B^2]\}}{dB} = \\ &= E\{-2y + 2Ax + 2B\} = 0 \end{aligned} \quad (5.154)$$

de unde rezultă mai departe că,

$$E\{-y + Ax\} + E\{B\} = 0 \Rightarrow B = E\{y - Ax\} = m_y - Am_x \quad (5.155)$$

Înlocuind în relația **Error! Not a valid link.** pe  $B$  cu valoarea sa dedusă în relația (5.219) obținem:

$$\begin{aligned} e &= E\{[y - (Ax + m_y - Am_x)]^2\} = E\{[(y - m_y) - A(x - m_x)]^2\} \Rightarrow \\ e &= E\{(y - m_y)^2\} - 2AE\{(y - m_y)(x - m_x)\} + A^2E\{(x - m_x)^2\} \Rightarrow \\ e &= \sigma_y^2 - 2A\rho\sigma_x\sigma_y + A^2\sigma_x^2 \end{aligned} \quad (5.156)$$

unde:  $\sigma_x^2$ ,  $\sigma_y^2$  sunt varianțele celor două variabile aleatoare,  $x$  și  $y$ , iar  $\rho$  este coeficientul de corelație Pearson (vezi **Subcapitolul 5.5.4**, relația (5.245)).

Mai departe relația (5.156), ce reprezintă eroarea  $e$  ca o funcție de  $A$ , este minimizată în raport cu  $A$  dacă:

$$\frac{de}{dA} = 0 \Leftrightarrow -2\rho\sigma_x\sigma_y + 2A\sigma_x^2 = 0 \Rightarrow A = \frac{\rho\sigma_y}{\sigma_x} \quad (5.157)$$

În concluzie, parametrii  $A$  și  $B$  pentru care se obține cel mai bun estimat liniar, în sensul erorii medii pătratice, al funcției de regresie a două variabile aleatoare  $x$  și  $y$  sunt dați de relațiile:

$$A = \frac{\rho\sigma_y}{\sigma_x} \text{ și } B = m_y - Am_x \quad (5.158)$$

**Observația 5.32:** În general, funcția de regresie neliniară nu este o linie dreaptă iar eroarea medie pătratică este mai mică decât eroarea în sens mediu pătratic din estimarea liniară regresiei.

**Problema 5.24:** Să se demonstreze că dacă variabilele aleatoare  $x$  și  $y$  au o densitate comună de tip *Gauss*-iană atunci, estimațiile, în sensul erorii medii pătratice minime, determinate prin cele două metode (liniară, respectiv, neliniară) pentru funcția de regresie sunt identice.

În cazul în care există mai mulți predictorii,  $x_i$ , în funcție de care se dorește exprimarea unei *v.a.* dependente  $y$ , atunci avem de-a face cu o regresie liniară multiplă iar modelul în acest caz este dat de:

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_k x_k + e \quad (5.159)$$

unde:  $x_i$ , cu  $i = \overline{1, k}$  sunt predictorii (regresorii) iar  $\alpha_i$  sunt coeficienții de regresie. Formula de determinare a coeficienților de regresie este una care diferă funcție de criteriul de eroare folosit. Cele mai multe însă din aplicațiile de estimare parametrică a regresiei liniare folosesc metoda celor mai mici pătrate.

În mod echivalent, ca în cazul estimării neliniare, și în cazul estimării liniare a regresiei numele metodei, respectiv, acela de ***estimare parametrică liniară*** se datorează faptului că forma funcției de regresie este una cunoscută, ea este dată de o combinație liniară a parametrilor modelului și, în plus, parametrii modelului au, așa cum am văzut, o semnificație particulară pentru procesul studiat.

**Observația 5.33:** În cazul particular discutat anterior am avut de-a face cu o regresie liniară simplă (o singură *v.a.* independentă) iar funcția de regresie estimată a fost o linie dreaptă. Nu trebuie să facem însă confuzie între numele metodei (estimare liniară) și această linie dreaptă care nu este decât un caz particular al funcției de regresie liniare. Mai mult chiar, inclusiv o funcție de forma:

$$y = \alpha_0 + \alpha_1 x_1^3 + \dots + \alpha_k x_k^2 + e \quad (5.160)$$

care este evident o funcție neliniară în variabilele independente  $x_i$ , este o funcție de regresie liniară întrucât, să ne reamintim – caracterul liniar al metodei este dat de liniaritatea exprimată în termenii parametrilor modelului – ceea ce se verifică și pentru relația noastră.

### **Estimarea parametrică a unui semnal**

Pe lângă aplicațiile de modelare a densităților de probabilitate sau a funcțiilor de regresie, estimarea parametrică își mai găsește aplicabilitate și în domenii cum ar fi modelarea unui semnal (pentru definiția unui semnal vezi **Anexa: Noțiunea de semnal. Semnal staționar**).



**Estimarea unui semnal** presupune găsirea unui model matematic care să fie, pe de o parte, cât mai apropiat de realitatea fizică care a generat semnalul observat iar pe de altă parte, același model trebuie să prezinte și un grad mare de generalitate.

Considerăm în continuare problema standard a unui semnal  $S$  peste care s-a suprapus un zgomot  $Z$ :

$$x_i = s_i + z_i, 1 \leq i \leq M \quad (5.161)$$

În ecuația de mai sus  $z_i$  sunt variabile aleatoare i.i.d., distribuite normal, după o distribuție  $N(0, \sigma^2)$  iar  $x_i$  sunt datele observate ale semnalului.

Estimarea, în acest caz, presupune că se cunoaște modelul care ar descrie procesul ce a generat datele semnalului și se dorește, pe baza datelor observate să se estimeze parametrii modelului, care descriu modelul, având totodată o anumită semnificație pentru acesta.

**Observația 5.34:** Un exemplu ilustrativ de estimare parametrică a unui semnal este și aplicația în care se urmărește modelarea unui semnal discret despre care se știe că este sinusoidal și pentru care se dispune de date observate. Precizăm aici că, o sinusoidă este orice funcție de forma  $A \cdot \sin(\omega t + \phi)$ , unde  $t$  este variabila independentă iar  $A$ ,  $\omega$  și  $\phi$  sunt parametrii fișii ai sinusoidei, numiți amplitudinea, frecvența și, respectiv, faza semnalului.

În acest caz, cel mai simplu model propus pentru semnalul  $s$  sinusoidal este dat de o singură componentă sinusoidală complexă,  $A \exp(j\omega n)$ . Rescriind relația (B036) pentru această aplicație particulară obținem:

$$x[n] = A e^{j\omega_0 n} + z[n] \quad (5.162)$$

cu  $1 \leq n \leq M$ , iar  $x[n]$  sunt datele observate ale semnalului, iar  $\theta = [A, \omega, \phi]^T$  este vectorul parametru al modelului propus, vector ce trebuie estimat. Estimarea parametrilor modelului matematic (parametri care, așa cum am văzut, caracterizează procesul analizat, având o semnificație particulară) presupune folosirea analizei statistice, respectiv, definirea unui criteriu de eroare. La ora actuală există propuse mai multe criterii de eroare însă dintre acestea, de departe cea mai folosită este metoda celor mai mici pătrate, în care se urmărește minimizarea sumei pătratelor erorii dintre datele observate și datele generate de model:

$$J(\theta) = \sum_{n=0}^{M-1} |x[n] - \hat{x}[n]|^2 \quad (5.163)$$

Minimizarea lui  $J(\theta)$  se face în raport cu vectorul parametru.

În estimarea parametrilor modelului unui semnal se pot folosi atât metode clasice de estimare statistică (de exemplu, metoda celor mai mici pătrate, metoda momentului, etc.), cât și metode *Bayes*-iene (de exemplu, metoda probabilității posterioare maxime), filtrul Kalman ș.a.m.d.

## 2. Estimarea neparametrică și semiparametrică

Așa după cum am amintit deja, este posibil să se facă inferențe statistice și fără să se presupună *a priori* o anumită familie parametrică a distribuțiilor de probabilitate sau un anumit model pentru funcția de regresie sau semnalul analizat. În acest caz vorbim de modele neparametrice (respectiv, de statistici neparametrice<sup>43</sup>) și de modele semiparametrice.

### Modelele neparametrice

**Modelele neparametrice** diferă de cele parametrice prin aceea că structura modelului este una mult mai flexibilă, ea nefiind specificată *a priori* ci ea este determinată direct din date. Din acest motiv, modelele neparametrice mai sunt numite și *modele ce nu depind de distribuție*.

Anticipând puțin, multe dintre metodele neparametrice folosesc în cursul procesului de estimare, așa după cum vom vedea, funcții de forma:

$$f(a) = \sum_{i=1}^{\infty} \alpha_i \varphi_i(a) \quad (5.164)$$

unde  $\alpha_i$  sunt considerați „parametrii” estimatorului neparametric iar  $\varphi_i(\cdot)$  reprezintă niște *funcții bază* (vezi **Anexa: Vectori bază și funcții bază**).

După cum se poate vedea termenul *neparametric*, asociat oarecum impropriu acestor modele statistice, nu semnifică faptul că astfel de modele sunt lipsite complet de parametri ci, faptul că:

- pe de o parte, numărul și natura parametrilor sunt flexibile și nu fixate dinainte (cu alte cuvinte, dacă cantitatea datelor crește este posibil să avem o funcție mai complexă și, deci, mai flexibilă), iar
- pe de altă parte, parametrii acestor modele nu au asociată o anumită semnificație, respectiv, aceea de parametri ai populației; în consecință, metodele neparametrice de estimare nu estimează parametrii populației.

<sup>43</sup> Funcții aplicate unui eșantion, a căror interpretare nu depinde de aproximarea vreunei distribuții de probabilitate parametrizate.

Prin estimare neparametrică înțelegem, în general, una din următoarele situații:

- *estimare de semnale* (vezi **Anexa: Noțiunea de semnal. Semnal staționar**),
- *estimare de funcții de regresie* și, respectiv,
- *estimare de distribuții de probabilitate*.

În mare vorbind, există două tipuri de metode de estimare neparametrică:

**Tabel 5.10.** Metode de estimare neparametrică

Metode de estimare neparametrică	Comentarii:
<b>Locală</b>	<ul style="list-style-type: none"> <li>- În această abordare estimarea unei funcții de valori reale<sup>44</sup>, <math>f(\cdot)</math>, se face punctual, în fiecare punct <math>a = a^0</math>. Din această perspectivă, estimarea funcției se reduce la estimarea succesivă a numerelor reale <math>f(a)</math>, unde variabila <math>a</math> baleiază tot domeniul de valori pe care este definită funcția (sau doar un domeniu de interes).</li> <li>- O altă estimare, tot punctuală, a funcției <math>f(\cdot)</math>, aplicabilă însă doar pentru cazul când aceasta este <u>o funcție continuu diferentiabilă</u><sup>45</sup>, este aceea care evaluează funcția nu doar în punctul generic <math>a^0</math> ci într-o vecinătate a acestuia, caz în care aproximarea punctuală a funcției se face prin valoarea <math>f(a^0) + f'(a^0)(a - a^0)</math>, ceea ce presupune, suplimentar, utilizarea observațiilor din vecinătate. La fel ca și mai sus, estimarea funcției presupune estimări punctuale succesive ce acoperă domeniul de interes.</li> <li>- Exemple de metode ce se bazează pe acest gen de abordare locală a estimării neparametrice sunt: histograma, estimarea cu funcții nucleu, metoda</li> </ul>

<sup>44</sup> Aici, prin funcție înțelegem oricare dintre următoarele variante: funcție densitate, funcție de regresie sau un semnal oarecare.

<sup>45</sup> Așa cum după cum știm, metodele neparametrice nu impun, în general, restricții asupra funcțiilor ce se estimează, cu excepția anumitor metode pentru care este necesară îndeplinirea de către funcția de estimat a condiției, nu foarte restrictive, de continuitate a funcției.

	celui mai apropiat eșantion ( $NN$ – Nearest Neighbour), metoda celor mai apropiate $k$ -eșantioane (metoda $k$ -NN) etc.
<b>Globală</b>	<p>- Aceste metode introduc un sistem de coordonate (în particular, o bază) într-un spațiu de funcții (vezi <b>Anexa: Vectori bază și funcții bază</b>); în acest mod problema estimării unei funcții se reduce la estimarea unui set de numere reale (este vorba de scalarii <math>\alpha_i</math> ce reprezintă coordonatele funcției în raport cu baza de funcții considerată).</p> <p>- Cu alte cuvinte, folosind o mulțime adecvată de funcții liniar independente <math>\{\varphi_i(a)\}_{i=1}^{\infty}</math> ca și coordonate ale spațiului de funcții atunci, orice funcție de valori reale și de energie finită poate fi, în mod unic, exprimată printr-o combinație liniară a funcțiilor bazei, respectiv, prin setul de coeficienți scalari <math>\alpha_i</math> astfel:</p> $f(a) = \sum_{i=1}^{\infty} \alpha_i \varphi_i(a) \quad (5.165)$ <p>- Alegerea bazei (care, așa cum am mai spus în Anexa.aaa, nu este unică într-un spațiu de funcții considerat) reflectă presupunerile noastre despre caracteristicile semnalului (de exemplu, putem avea de-a face cu un semnal lent variabil, continuu, semnal oscilator ș.a.m.d.).</p> <ul style="list-style-type: none"> <li>• niște baze bine cunoscute sunt <i>seriile polinomiale</i> și <i>seriile Fourier</i> (aceste funcții bază au proprietatea că sunt infinit diferentiabile în orice punct);</li> <li>• alte funcții larg folosite sunt și funcțiile bază spline polinomiale și funcțiile bază wavelet etc.</li> </ul>

Pentru a înțelege următorul exemplu pe care îl vom discuta, facem aici precizarea că funcțiile nucleu folosite în estimarea neparametrică de tip local permit, prin însăși modul cum sunt definite (vezi **Subcapitolul 8.2.2**), o „analiză” a vecinătății fiecărui punct în care are loc aproximarea funcției dorite. Această analiză presupune de fapt asignarea de ponderi fiecărei

observații particulare aflate în vecinătatea punctului în care se face estimarea; în acest caz, vecinătatea este definită de un parametru al funcției nucleu, numit lățime de bandă sau parametrul de netezire al funcției iar ponderile sunt o măsură a proximității acestor observații față de punctul ales, ele fiind întotdeauna zero pentru observațiile din afara vecinătății.

**Exemplu 5.9:** Problema estimării neparametrice a unei funcții de regresie (cazul univariat) folosind pentru aceasta funcțiile nucleu, presupune că, în formula dată de relația (5.215):

$$b = m(a; \theta) + e,$$

pentru  $m(a; \theta)$  – reprezentând media condiționată a v.a.  $y$  ( $b$ , variabilă generică) dată de v.a.  $x = a$  –, nu se precizează *a priori* nici o formă parametrică. Cele  $N$  perechi de observații  $(a_i, b_i)$  ale eșantionului sunt folosite pentru a estima funcția densitate comună pentru variabilele aleatoare  $x$  și  $y$  astfel: densitatea în punctul  $(a^0, b^0)$  este estimată prin observarea, din cele  $N$  date empirice, a proporției de date care se află “aproape” de punctul  $(a^0, b^0)$ . Această procedură implică utilizarea unei funcții nucleu ce asignează ponderi corespunzătoare observațiilor din vecinătatea punctului în care se face estimarea.

Modelele neparametrice, deși prezintă avantajul unei foarte mari flexibilități, ele au și o serie de dezavantaje:

- (a) Pentru a aproxima, în mod adecvat, o funcție, modelul folosit (vezi, de exemplu, relația (5.164)) are un număr practic infinit de parametri necesari pentru a putea “captura” caracteristicile complexe, presupuse, ale funcției de estimat. Un număr mare al parametrilor modelului (parametri ce trebuie estimați) duce la o rată de convergență a metodei foarte mică. Dependența ratei de convergență a unui estimator de numărul parametrilor modelului este adesea referită sub numele de “blestemul dimensionalității”.
- (b) Pe de altă parte, pentru a estima într-un mod corespunzător parametrii de mai sus avem nevoie de un eșantion de date suficient de mare (condiție ce în practică arareori poate fi îndeplinită).
- (c) Un alt dezavantaj constă în aceea că, în mod deosebit pe seturi mari de date, acești estimatori flexibili presupun un volum de calcul foarte mare.

### **Modele semiparametrice**

Modelarea semiparametrică a apărut ca un compromis între modelarea parametrică și cea neparametrică, compromis concretizat prin introducerea

parțială în model a unor componente parametrice. Din acest punct de vedere, modelele semiparametrice sunt mai restrictive decât modelele complet nespecificate dar mai flexibile însă decât modelele parametrice. Modelele semiparametrice sunt construite, în general, pornind de la modele neparametrice în care s-au introdus, parțial, componente parametrice (parametri ce au o anumită semnificație pentru populația sau procesul analizat).

Metodele semiparametrice prezintă, ca un avantaj major, rate de convergență mult mai mari decât ratele de convergență obținute în abordările neparametrice ale aceluiași probleme. Mai mult chiar, în multe cazuri rata de convergență obținută este una similară cu cea obținută folosind metode parametrice. Câteva exemple de modele semiparametrice sunt și: regresia generalizată, modelele liniare parțial generalizate, modelele aditive și modelele aditive generalizate etc.

### **Histograma**

În practică, ori de câte ori dispunem de un set de date empirice pentru un vector aleator ne este foarte dificil să tragem concluzii cu privire la proprietățile acestuia din urmă – și aceasta folosindu-ne doar de simpla inspecție a valorilor observate. Într-un astfel de moment ne vine în ajutor statistica descriptivă care, fie prin *reprezentări grafice*, fie prin *măsurile cantitative*, extrage informații asupra caracteristicilor datelor la nivel de eșantion, informații pe care le putem utiliza ulterior pentru a face inferențe la nivel de populație sau pentru a lansa anumite ipoteze de lucru.

**Exemplu 5.10:** O astfel de ipoteză de lucru este, spre exemplu, și cea folosită în cadrul procesului de estimare parametrică în care forma funcțională a densității de probabilitate a caracteristicii la nivel de populație, deși necunoscută, se presupune a fi una cunoscută. În acest caz, la baza alegerii *a priori* a uneia sau alteia dintre densitățile teoretice existente stau deopotrivă: **(a)** o cunoaștere prealabilă a familiilor de distribuții teoretice și **(b)** informațiile culese fie printr-o analiză științifică a fenomenului fizic studiat, fie printr-o analiză empirică a datelor observate (vorbim aici deci, de statistica descriptivă aplicată la nivel de eșantion).

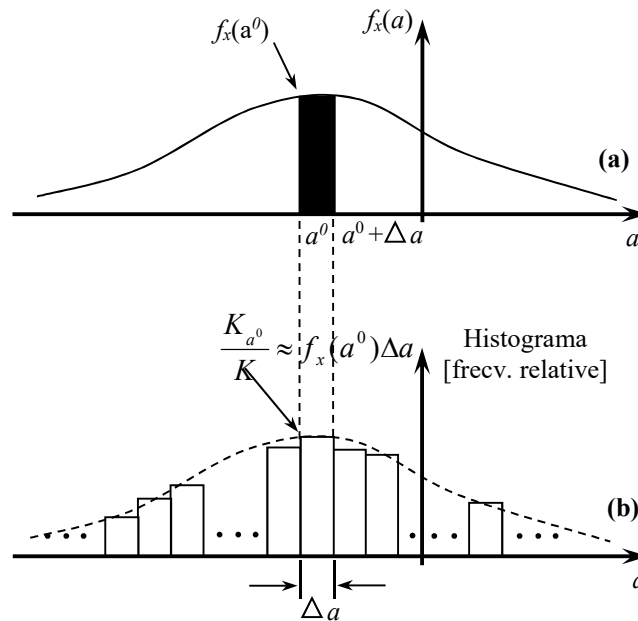
Cea mai directă modalitate de investigare a caracteristicilor unui eșantion este dată de prezentarea, într-o formă grafică convenabilă, a datelor observate. O astfel de afișare grafică – larg utilizată doar pentru *datele univariate și bivariate* – este și histograma.

**Histograma** este, prin definiție, o reprezentare grafică, sub formă de bare, a frecvențelor relative sau absolute de apariție a unei variabile aleatoare, pe intervale de valori. Mai exact, ea ilustrează proporția realizărilor particulare ale  $v.a.$   $x$  care se regăsesc în fiecare dintre cele câteva intervale disjuncte de valori, predefinite, ale lui  $x$ . Aceste intervale de valori sunt adiacente și, în mod ideal, de lățime egală.

Dacă lățimea unui interval de valori tinde, la limită, către o singură valoare iar numărul de intervale (respectiv, puncte) pe care se face reprezentarea histogramei tinde la infinit atunci histograma devine, așa cum vom vedea mai jos, chiar funcția densitate de probabilitate ce caracterizează setul de date empirice.

În realitate, nu vom dispune niciodată de un număr infinit de date empirice și, în consecință, histograma, așa cum a fost ea definită mai sus doar *aproximează funcția densitate de probabilitate* asociată, fiind, prin modul de abordare, un *estimator neparametric* al acesteia. Estimatul obținut pe un eșantion poartă, în acest caz, numele de *distribuție empirică*.

În **Figura 5.20(a)** este reprezentată grafic o posibilă formă a unei funcții densitate de probabilitate unidimensională  $f_x(a)$ , necunoscută, ce este modelată prin histograma din figura **Figura 5.20(b)**.



**Figura 5.20.** (a). Funcția densitate de probabilitate a unei variabile aleatoare  $x$  versus (b) histograma normalizată corespunzătoare

Histograma setului de date din **Figura 5.20(b)** a fost obținută printr-o segmentare pe lățimi  $\Delta a$  a intervalului de valori ale lui  $x$ . Înălțimea fiecărui dreptunghi este una proporțională cu numărul de elemente din respectivul interval de lățime  $\Delta a$ . Fie  $K_{a^0}$  numărul de valori ale v.a.  $x$  aflate în intervalul  $(a^0, a^0 + \Delta a]$ , respectiv  $a^0 < a \leq a^0 + \Delta a$ . În cazul general, când numărul total de date empirice,  $K$ , din eșantion este suficient de mare iar lățimea  $\Delta a$  a intervalelor este suficient de îngustă, ambele cantități  $K_{a^0}/K$  cât și  $f_x(a^0)\Delta a$  sunt, de fapt, aceleași estimări ale probabilității de poziționare a unei realizări particulare a v.a.  $x$  în intervalul anterior prezentat (vezi **Anexa: Media aritmetică(a)**). Astfel:

$$f_x(a^0)\Delta a \approx \frac{K_{a^0}}{K} \quad (5.166)$$

În esență vorbim aici de *interpretarea frecvențială* a probabilității (pentru alte interpretări a se vedea **Anexa: Interpretarea noțiunii de probabilitate**).

O discuție ce s-ar impune a fi făcută aici ține de alegerea optimă a numărului de intervale de valori,

$$N_{\Delta a} = \text{int}\left(\frac{\max(a) - \min(a)}{\Delta a}\right) \quad (5.167)$$

pe care se reprezintă histograma<sup>46</sup> (aici  $\text{int}(\cdot)$  reprezintă partea întregă). Lățimi diferite ale acestor intervale pun, în general, în evidență caracteristici diferite ale datelor. Nu există la ora actuală metode care să ne ofere o soluție optimă la această problemă și, de aceea, alegerea finală se va face prin încercări succesive; respectiv, se dau pe rând valori (nu foarte mari) lui  $\Delta a$  și se reține cea valoare care pune în evidență anumite particularități ascunse ale datelor. Dintre formulele propuse pentru calculul acestui parametru important al histogramei prezentăm aici, ca informație orientativă în alegerea inițială a lui  $\Delta a$ , formula lui Scott:

$$\Delta a = 3.5 \cdot \sigma \cdot K^{-1/3} \quad (5.168)$$

<sup>46</sup> O astfel de discuție își are sensul dacă avem mereu reprezentat în minte faptul că histograma se realizează, în general, pe un eșantion de date ales în mod aleator din populația de interes. În consecință, frecvențele de apariție în cadrul eșantionului ale anumitor valori s-ar putea să difere uneori considerabil de mult față de frecvențele de apariție în populație ale aceluiași valori; concret, vorbim de valori cu mult mai mici ale acestor frecvențe. În acest context, o histogramă realizată pe intervale foarte mici s-ar putea să genereze o reprezentare foarte zgomotoasă care, practic, nu ne furnizează foarte multă informație privind distribuția reală a datelor.



unde  $\sigma$  reprezintă deviația standard a v.a.  $x$ .

Informațiile pe care ni le furnizează histograma sunt informații utile privind caracteristicile datelor, cum ar fi **tendința centrală**, **dispersia datelor**, precum și **forma (alura) generală a distribuției**. Vom vedea în subcapitolul următor ce importanță majoră poate juca histograma în inferența statistică, în special prin această ultimă informație furnizată.

**Observația 5.35:** Revenind în contextul **Exemplu 5.10**, o inspecție vizuală a histogramei eșantionului ne poate sugera, de exemplu, o anumită formă funcțională pentru densitatea de probabilitate necunoscută la nivel de populație, formă pe care o folosim, mai departe, ca *ipoteză de lucru* în estimarea parametrică. Ulterior fixării acestei ipoteze de lucru, în continuare – în cadrul procesului de inferență – se vor mai estima doar parametrii asociați respectivei funcții densitate de probabilitate.

**Observația 5.36:** Așa după cum am prezentat la început, histograma este, în principal, un instrument grafic de vizualizare pentru domeniul univariat și bivariat de valori. Această caracteristică nu trebuie privită ca fiind una limitativă în ceea ce privește utilitatea metodei în cazul vectorilor aleatori multidimensionali. Și în acest caz, o inspecție vizuală pe fiecare variabilă aleatoare componentă, combinată eventual cu alte informații statistice (de exemplu, de independență statistică între componentele vectorului etc.), ne poate conduce la niște concluzii la fel de utile privind forma funcțională a fdp caracteristică vectorului.

**Observație 5.37:** Așa după cum am prezentat la început, histograma este, în principal, un instrument grafic de vizualizare pentru domeniul univariat de valori însă există și histograme multivariate.

În prezent, în neurobiologie, ca și în multiple alte domenii, histograma este în continuare folosită cu succes ca un instrument tradițional de estimare a funcției densitate de probabilitate a procesului aleator studiat, deși alte metode neparametrice de estimare a densității s-au dovedit a fi mai eficiente. Această situație se explică în principal prin simplitatea metodei și prin informația mai intuitivă, de natură vizuală, pe care ne-o furnizează.

### 5.5.2. Momentele unui vector aleator

Pe lângă tehnicile de vizualizare grafică a trăsăturilor unui set de date empirice, statistica descriptivă oferă și o serie de *măsuri cantitative* ce

sumarizează, într-o manieră exactă, caracteristicile respectivului eșantion de date.

Astfel, informațiile pe care ni le furnizează histograma pot fi extrase și exprimate și în termeni pur cantitativi, prin următoarele măsuri:

- tendința centrală, prin *medie* și *mediană*, respectiv,
- dispersia datelor, prin *rang* și *deviația standard*.

Parametrii<sup>47</sup> tendinței centrale reprezintă cele mai importante măsuri folosite în caracterizarea unei distribuții empirice. Valorile furnizate de acestea ne ajută să localizăm datele pe o scală liniară, oferindu-ne totodată o imagine asupra celei mai bune valori ce descrie datele. Dintre acești parametri cei mai uzuali sunt *media aritmetică* a eșantionului (vezi relația **Error! Not a valid link.**) și, respectiv, *mediana* – aceasta reprezentând acea valoare a vectorului aleator  $x$  pentru care 50% din datele eșantionului au valori sub această valoare iar 50% au valori peste această valoare. Mediana este folosită, de obicei, ca alternativă a mediei aritmetice care este un parametru sensibil la valorile eronate ale eșantionului (valori extreme care diferă foarte mult de majoritatea datelor). Acest ultim neajuns nu constituie însă o problemă în sine, întrucât, în general, respectivele valori eronate sunt eliminate într-o așa-zisă fază de preprocesare a datelor.

Dintre parametrii ce măsoară gradul de împrăștiere al datelor am amintit mai sus *rangul*, calculat ca diferență între valoarea maximă și valoarea minimă din setul de date (parametru extrem de sensibil în cazul unor date extreme, eronate), respectiv *deviația standard* calculată ca media abaterii fiecărui punct de la media calculată a eșantionului. O altă măsură a dispersiei datelor este și *varianța*, care reprezintă pătratul deviației standard.

Așa după cum am văzut în **Subcapitolul 5.5.1**, subpunctul 1 „**Estimarea parametrică punctuală clasică**”, atât media cât și varianța distribuției empirice sunt parametri de selecție și ei sunt adesea utilizați ca estimări ale parametrilor statistici corespunzători, din populație.

După cum vom vedea în cele ce urmează, atât media cât și varianța reprezintă, în fapt, cele mai utilizate două măsuri cantitative asociate unei v.a. dintr-o, practic, infinitate de măsuri definite teoretic și care sunt cunoscute generic sub numele de **momente**.

<sup>47</sup> Se cunosc mai multe măsuri cantitative ce descriu tendința centrală. Dintre acestea amintim aici și *modul*, *media geometrică*, *media armonică* etc.

## 1. Momentele în analiza statistică. Funcții caracteristice

**Momentele** unei distribuții reprezintă cantități de interes în studiul variabilelor și vectorilor aleatori. Mai jos avem definite **momentele pentru o v.a.  $x$** :

Momentele: 
$$m_n = E\{x^n\} = \int_{-\infty}^{+\infty} a^n f_x(a) da$$

Momentele centrate: 
$$\mu_n = E\{(x - m_x)^n\} = \int_{-\infty}^{+\infty} (a - m_x)^n f_x(a) da$$

Momentele absolute: 
$$E\{|x|^n\}, \quad E\{|x - m_x|^n\}$$

Momentele generalizate: 
$$E\{(x - \alpha)^n\}, \quad E\{|x - \alpha|^n\}$$

În practică, cele mai utilizate sunt momentele și momentele centrate de ordin 1 și, respectiv, 2:  $m_1 = m_x, \mu_2 = \sigma_x^2$ .

Trebuie să reținem faptul că momentele unei v.a. nu sunt numere arbitrare ci ele satisfac numeroase inegalități cum ar fi, de exemplu:  $\sigma^2 = m_2 - m_1^2 \geq 0$  (vezi **Problema 5.26** rezolvată în anexă).

Pe lângă momentele prezentate mai sus există definite și momentele centrate normalizate (standardizate); acestea sunt date de momentul centrat de ordin  $n$  corespunzător, divizat prin  $\sigma^n$ , respectiv,

$$\text{moment centrat de ordin } n \text{ normalizat} = \mu_n / \sigma^n \quad (5.169)$$

Aceste momente normalizate sunt cantități fără unitate de măsură ce reprezintă distribuția într-un mod independent de orice modificare liniară a scalei.

În mod corespunzător, **momentele comune pentru două variabile aleatoare**,  $x$  și  $y$ , se definesc și ele astfel:

Momentele: 
$$m_{kl} = E\{x^k y^l\} = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} a^k b^l f_{xy}(a, b) dadb$$

Media produsului  $x^k y^l$  reprezintă momentul comun de ordin  $n = k + l$  al variabilelor aleatoare  $x$  și  $y$ .

Momentele centrate: 
$$\begin{aligned} \mu_{kl} &= E\{(x - m_x)^k (y - m_y)^l\} = \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (a - m_x)^k (b - m_y)^l f_{xy}(a, b) dadb \end{aligned}$$

Momentele absolute: 
$$E\{|x|^k \cdot |y|^l\}, \quad E\{|x - m_x|^k \cdot |y - m_y|^l\}$$

Momentele generalizate:  $E\{(x-\alpha)^k(y-\beta)^l\}$ ,  $E\{|x-\alpha|^k \cdot |y-\beta|^l\}$

Din definițiile de mai sus rezultă că:

- *momentele comune de ordin unu* sunt:  $m_{10} = m_x$  și  $m_{01} = m_y$ ;
- *momentele comune de ordin doi* sunt:  $m_{20} = E\{x^2\}$ ,  $m_{11} = E\{xy\}$  și  $m_{02} = E\{y^2\}$ ;
- *momentele centrate comune de ordin unu* sunt:  $\mu_{10} = \mu_{01} = 0$ ;
- *momentele centrate comune de ordin doi* sunt:  $\mu_{21} = \sigma_x^2$ ,  $\mu_{11} = Cov(x,y)$ <sup>48</sup> și  $\mu_{02} = \sigma_y^2$ .

**Observația 5.38:** Definițiile momentelor (respectiv, a momentelor comune) pentru vectori aleatori sunt similare celor prezentate, mai sus, pentru variabile aleatoare, cu particularitatea, bineînțeles, a folosirii operatorului de transpunere în cazul produsului scalar al vectorilor.

Utilitatea momentelor este una multiplă:

- astfel, am văzut până acum rolul jucat de acestea în estimarea parametrică (statistica inferențială) – acolo unde anumite momente coincid chiar cu parametrii distribuției teoretice ce se presupune a modela funcția de densitate necunoscută a unei populații; în acest caz, doar calcularea respectivelor momente la nivel de eșantion conduc direct la obținerea unui estimat al densității de probabilitate dorite;
- o altă utilitate a momentelor definite mai sus, întâlnită de această dată în statistica descriptivă, rezidă în capacitatea lor de a furniza o serie de informații descriptive privind datele analizate; în cele ce urmează vom discuta această ultimă direcție de utilizare a momentelor.

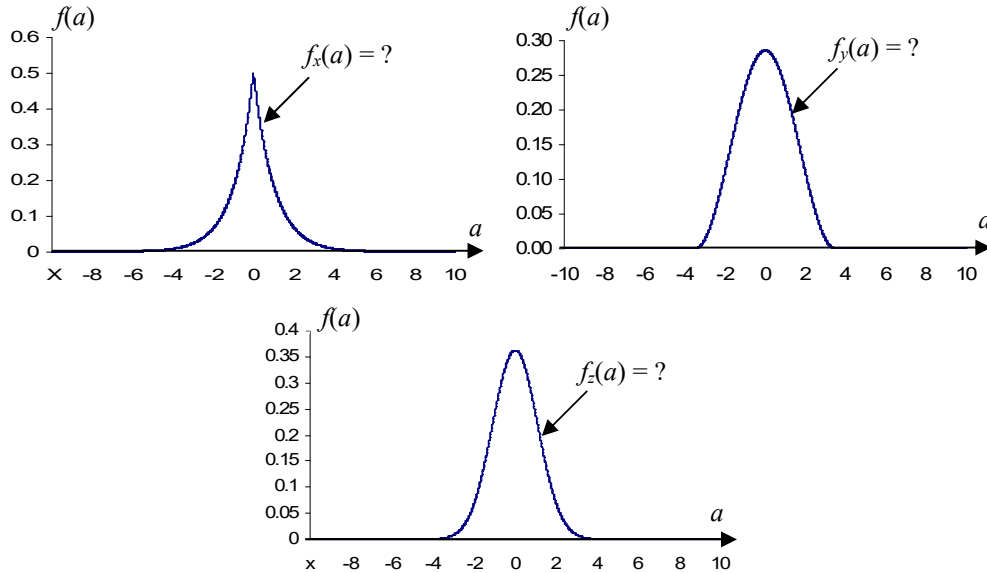
Un rezultat important legat de momentele unei v.a. (respectiv, vector aleator)  $x$  este cel al unei teoreme, numită **teorema momentului**, care afirmă că: *dacă  $m_n$  este cunoscut pentru fiecare  $n$ , atunci, în anumite condiții* (vezi mai jos), *funcția densitate de probabilitate a lui  $x$  este determinată în mod unic*. O prezentare pe scurt a acestei teoreme, precum și a noțiunilor noi pe care le utilizează aceasta, o redăm în tabelul de mai jos:

Funcții ale unei v.a. $x$	Definiții/Enunț
<i>Funcția caracteristică</i>	Notată $\Phi_x(\omega)$ , este prin definiție dată de integrala, $\Phi_x(\omega) = \int_{-\infty}^{+\infty} f_x(a)e^{j\omega a} da = E\{e^{j\omega x}\} \quad (5.170)$

<sup>48</sup> Coeficientul de covarianță al variabilelor aleatoare  $x$  și  $y$ .

<p><b>Funcția (generatoare) de momente</b></p>	<p>Notată, <math>M_x(s)</math>, se obține din funcția caracteristică astfel: pentru <math>M_x(j\omega) = \Phi_x(\omega)</math> și <math>s = j\omega</math>,</p> $M_x(s) = \int_{-\infty}^{+\infty} f_x(a) e^{sa} da = E\{e^{sx}\} \quad (5.171)$
<p><b>Teorema momentului</b></p>	<p>Derivând (5.235) de <math>n</math> ori obținem,</p> $M^{(n)}(s) = E\{x^n e^{sx}\} \quad (5.172)$ <p>și, corespunzător,</p> $M^{(n)}(0) = E\{x^n\} = m_n \quad (5.173)$ <p>Dacă toate momentele <math>m_n</math> sunt finite și dacă, în plus, dezvoltarea în serie a lui <math>M_n(s)</math> în apropierea originii,</p> $M_x(s) = \sum_{n=0}^{\infty} \frac{m_n}{n!} s^n \quad (5.174)$ <p>este o serie convergentă, atunci <i>densitatea</i> <math>f_x(a)</math> – ce poate fi scrisă în termenii lui <math>M_n(s)</math> (vezi și relația (5.239)) – este, atunci când se cunosc toți <math>m_n</math>, <i>unic determinată</i>.</p>
<p><i>Notă:</i></p> <p>(a) Întrucât <u>funcția caracteristică</u>, <math>\Phi_x(\omega)</math>, este strâns legată de transformata Fourier – mai exact, <math>\Phi_x(\omega)</math> a lui <math>x</math> este <u>complex-conjugatul</u> transformatei Fourier continue a funcției de densitate a lui <math>x</math> (vezi <b>Anexa: Transformata Fourier</b>), atunci proprietățile lui <math>\Phi_x(\omega)</math> sunt, în esență, aceleași cu proprietățile transformatelor Fourier. În acest caz, <math>f_x(a)</math> poate fi exprimat în termenii lui <math>\Phi_x(\omega)</math>, astfel:</p> $f_x(a) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} \Phi(\omega) e^{-j\omega a} d\omega \quad (5.175)$ <p>(b) Dacă <math>x = [x_1, \dots, x_d]^T</math> este un vector aleator <math>d</math>-dimensional atunci argumentul scalar <math>\omega</math> devine vectorul linie, <math>d</math>-dimensional, <math>\Omega = [\omega_1, \dots, \omega_d]</math> iar produsul <math>\omega x</math> devine produsul scalar <math>\Omega x</math>:</p> $\Phi_x(\Omega) = E\{e^{j\Omega x}\} = E\{e^{j(\omega_1 x_1 + \dots + \omega_d x_d)}\} = M_x(j\Omega) \quad (5.176)$ <p>(c) Funcția caracteristică a unei v.a. există întotdeauna ceea ce nu este și cazul funcției generatoare de moment. În cazul în care aceasta din urmă nu există funcția caracteristică este cea utilizată pentru a determina momentele.</p> <p>În concluzie, <i>orice distribuție de probabilitate (pe <math>R</math> sau pe <math>R^n</math>) are o funcție caracteristică și, reciproc, pentru orice funcție caracteristică există exact o singură distribuție de probabilitate</i>. Altfel spus, funcția caracteristică a unei v.a. <i>definește complet și în mod unic distribuția de probabilitate</i> a respectivei variabile.</p>	

O concluzie imediată ce se poate trage de aici este și faptul că, momentele de ordin unu și doi – de departe cele mai utilizate în aplicațiile practice – nu dau decât o caracterizare limitată a funcției de densitate necunoscută a populației. În acest context, cunoașterea și a altor momente poate fi utilă mai ales în faza de alegere între două densități ce prezintă aceleași momente de ordin unu și doi și care sunt candidate la calitatea de estimat al unei funcții de densitate de probabilitate necunoscute; în consecință, momentele ne pot asista în mod efectiv în procesul de selecție a acelei distribuții teoretice care se apropie cel mai mult de alura (dezvăluită prin indicatori sintetici adecvați) a funcției de densitate reale.



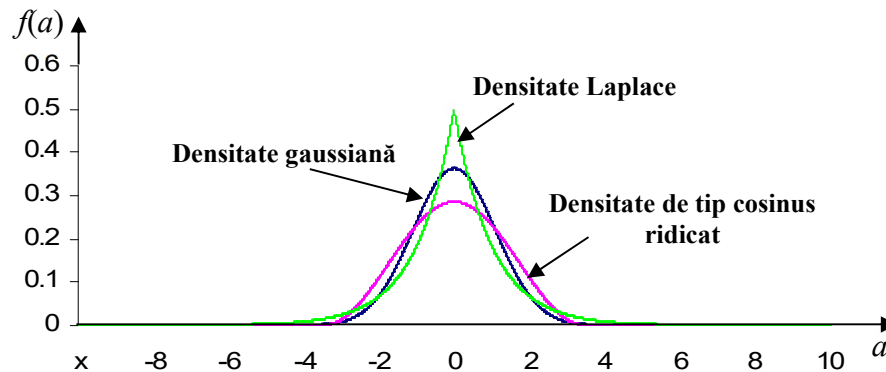
**Figura 5.21.** Funcții de densitate de probabilitate de aceeași medie și varianță

**Exemplu 5.11:** În **Figura 5.21** sunt reprezentate trei densități empirice unidimensionale (obținute cu ajutorul histogramei) – densități care, prezintă aceeași medie,  $m$ , și aceeași varianță,  $\sigma^2$ . Sunt aceste densități din familia de distribuții gaussiene?

În urma unei prime inspecții vizuale am fi tentați să spunem că din cele trei densități numai ultimele două ar corespunde descrierii așa numitului „clopot al lui Gauss” (vezi **Subcapitolul 5.6.2**).

O suprapunere a celor trei grafice (vezi **Figura 5.22**) ne dezvăluie însă o situație destul de contradictorie și anume, suntem în situația în care pentru aceiași doi parametri, *medie* și, respectiv, *varianță* –

parametri ce știm că definesc complet o distribuție gaussiană –, obținem densități *Gauss*-iene distincte.



**Figura 5.22.** Forme funcționale diferite pentru trei densități cu aceeași medie și aceeași varianță, însă cu coeficient de boltire diferit.

Răspunsul la această dilemă îl vom afla, însă, abia în momentul în care vom calcula valoarea unor **indicatori sintetici** suplimentari (momente centrate normalizate, vezi (5.169)), indicatori care ne oferă informații privitoare de această dată la **alura funcției** de densitate empirică.

În principal, vorbim de două categorii de astfel de indicatori sintetici, respectiv:

- (a) indicatori de asimetrie, respectiv,
- (b) indicatori de boltire (sau ai excesului).

**Indicatorii de asimetrie** dau informații asupra modului de repartizare a frecvențelor de o parte sau alta a valorii centrale a unei serii aleatoare. În literatura de specialitate s-au propus mai mulți indicatori pentru măsurarea gradului de asimetrie al distribuțiilor. Cel folosit de noi în **Exemplul 5.11** a fost calculat cu formula:

$$\gamma_1 = \frac{\mu_3}{\sigma^3} \quad (5.177)$$

Valori ale lui  $\gamma_1$ :

- egale cu zero → indică o distribuție simetrică;
- strict negative → indică o distribuție mai împrăștiată spre stânga, față de valoarea medie, cu o coadă a distribuției mai mare spre stânga (*asimetrie de stânga*);

- strict pozitive → indică o distribuție mai împrăștiată spre dreapta, față de valoarea medie, cu o coadă a distribuției mai mare spre dreapta (*asimetrie de dreapta*).

**Indicatorii de boltire** dau informații asupra gradului cât sunt de ascuțite, respectiv, de plate, densitățile empirice față de o densitate de probabilitate normală. Dintre indicatorii propuși, prezentăm doar *coeficientul de boltire al lui Fisher*, respectiv:

$$\gamma_2 = \frac{\mu_4}{\sigma^4} \quad (5.178)$$

Valori ale lui  $\gamma_2$ :

- egale cu 3 → indică o distribuție normală (*distribuție mezokurtică*);
- mai mari decât 3 → indică o distribuție cu un vârf aproape de medie mai ascuțit decât cel al unei distribuții normale (*distribuție leptokurtică*);
- mai mici decât 3 → indică o distribuție cu un vârf aproape de medie mai aplatizat față de cel al unei distribuții normale (*distribuție platikurtică*).

**Observația 5.39:** Pentru o corectă interpretare a valorilor calculate, în momentul implementării sau a folosirii pentru indicatorii de mai sus a unor funcții gata predefinite în anumite medii de lucru (de exemplu, Matlab, Microsoft Office Excel etc.) trebuie avut grijă la următoarele aspecte: (a) dacă în relația (5.178) este sau nu substrasă valoarea 3 (în prima situație coeficientul de boltire calculat pentru o distribuție normală va avea valoarea zero în loc de trei) și (b) întrucât acești indicatori sunt calculați pe un eșantion iar nu pe întreaga populație, valorile obținute vor fi deplasate (vor tinde să difere de valorile lor reale din populație cu o cantitate sistematică ce depinde de mărimea eșantionului); în acest caz este indicată realizarea corespunzătoare a unei corecții a valorilor obținute.

Revenind la **Exemplul 5.11**, calculul coeficientului de asimetrie, respectiv, al coeficientului de boltire pentru cele trei densități prezentate ne conduce la următoarele rezultate – rezultate ce concordă cu caracteristicile (în particular, alura) funcțiilor teoretice folosite la generarea densităților empirice respective:



Densitate empirică	Momente normalizate	Distribuții teoretice folosite
$f_x(a)$	$\gamma_1 = 0, \gamma_2 = 6$	Distribuția Laplace: $f_x(a) = \frac{1}{2\beta} \exp\left(-\frac{ a - m_x }{\beta}\right) \quad (5.179)$ Calculată pentru $\beta$ (factorul de scalare) = 1.
$f_y(a)$	$\gamma_1 = 0, \gamma_2 = 2.7$	Distribuția normală, $N(0, 1.1)$ .
$f_z(a)$	$\gamma_1 = 0, \gamma_2 = 1$	Distribuția cosinus ridicat: $f_x(a) = \begin{cases} \frac{1}{2s} \left[ 1 + \cos\left(\frac{a - m_x}{s} \pi\right) \right] & \text{pentru }  a  \leq s \\ 0 & \text{în rest.} \end{cases} \quad (5.180)$ calculată pentru $s$ (ce definește limitele intervalului pe care funcția ia valori diferite de zero) = 3.5.

Notă: Implementarea coeficienților s-a făcut în Microsoft Excel™ iar dimensiunea eșantionului de date a fost  $n = 2000$ .

În concluzie, complexitatea realității impune folosirea unui sistem de indicatori. Altfel spus, atunci când alegem o distribuție teoretică sau alta pentru a modela densitatea necunoscută din populație, pe lângă informațiile furnizate de momentele de ordin unu și doi, de mare utilitate ne sunt și informațiile furnizate de momentele de ordin superior, ce ne dau o imagine mai corectă privind forma densității empirice. Din acest punct de vedere și histograma este o tehnică grafică eficientă ea punând în evidență această caracteristică a datelor.

O altă utilitate a momentelor, ce se încadrează tot în statistica descriptivă, este dată de informațiile descriptive pe care ni le furnizează, în principal, momentele comune de ordin unu și doi a două variabile aleatoare, respectiv, a doi vectori aleatori. Aceste informații utile (ca de exemplu, independența sau gradul de corelare a două au mai multe variabile aleatoare) rezidă în însăși proprietățile acestor momente care se definesc așa cum s-a prezentat mai sus.

În continuarea acestui subcapitol vom defini momentele de ordin I și II pentru un vector aleator, vom prezenta proprietățile acestora precum și câteva proceduri simple de estimare a acestor momente din setul de date. Reamintim că prin estimare parametrică se înțelege o *evaluare aproximativă* a parametrilor de interes.

## 2. Momentul de ordin întâi

Fie  $x$  un vector aleator, discret sau continuu. Atunci, *valoarea sperată* sau *media* (momentul statistic de ordinul I, notat cu  $m_x$ ) al acestui vector aleator este, prin *definiție*, dată de relația:

$$E\{x\} \stackrel{\text{def}}{=} \int_{-\infty}^{+\infty} a f_x(a) da \quad (5.181)$$

$$E\{x\} \stackrel{\text{def}}{=} \sum_{-\infty}^{+\infty} a p_x(a) \quad (5.182)$$

Relația (5.181) este pentru un vector aleator caracterizat de variabile aleatoare continue, în timp ce relația (5.182) este pentru un vector aleator caracterizat de variabile aleatoare discrete.

Dacă integrala sau suma date de relațiile anterioare, pentru cel puțin una dintre componentele vectorului aleator  $x$ , nu converge către un număr finit vom spune că estimatorul vectorului aleator  $x$  nu există.

Media unui vector aleator este un *operator liniar*, el fiind caracterizat de următoarele proprietăți.

### Proprietăți:

1. Pentru orice vectori aleatori  $x$  și  $y$  avem:

$$E\{x + y\} = E\{x\} + E\{y\} \quad (5.183)$$

2. Pentru orice vector aleator  $x$  și pentru orice constantă scalară  $\alpha$  și vector constant  $a^0$  avem:

$$E\{\alpha x + a^0\} = \alpha E\{x\} + a^0 \quad (5.184)$$

3. Fie  $x$  un vector aleator și fie  $g(x)$  o cantitate derivată din el. De exemplu, această cantitate,  $g(x)$ , poate fi scalară, vectorială sau matricială, precum:  $g(x) = x x^{*T}$  sau  $g(x) = (x - m_x)(x - m_x)^{*T}$ . Media (momentul statistic de ordinul I) a lui  $g(\cdot)$  este notată  $E\{g(x)\}$  și ea este definită de relațiile matematice:

$$E\{g(x)\} = \int_{-\infty}^{+\infty} g(a) f_x(a) da \quad (5.185)$$

$$E\{g(x)\} = \sum_{-\infty}^{+\infty} g(a) p_x(a) \quad (5.186)$$

Relația (5.185) este pentru un vector aleator caracterizat de variabile aleatoare continue în timp ce relația (5.186) este pentru un vector aleator ce conține variabile aleatoare discrete.

Dacă  $g(a)$  este un vector sau o matrice, în relația (5.185), se aplică estimarea statistică pentru fiecare element în parte iar rezultatul estimării statistice va fi un alt vector sau o altă matrice cu aceleași dimensiuni.

**Exemplu 5.12:** Pentru o posibilă aplicație a semnalelor aleatoare în analiza convertoarelor AD (analog – digitale) consultați **Anexa: Analiza raportul semnal/zgomot pentru un convertor ideal.**

**Media condiționată** a unui vector aleator  $x$ , presupunând evenimentul  $B$ , este dată de integrala (5.181) în care fdp  $f_x(a)$  este înlocuită de fdp condiționată  $f_x(a|B)$ :

$$E\{x|B\} = \int_{-\infty}^{+\infty} a f_x(a|B) da \quad (5.187)$$

Cele mai comune momente sunt cele de ordin unu și doi.

Momentul de ordin întâi sau media vectorului aleator este definit de relația:

$$m_x = E\{x\} \stackrel{def.}{=} \int_{-\infty}^{+\infty} a f_x(a) da \quad (5.188)$$

Deoarece  $x$  este un vector  $d$  dimensional,  $m_x$  este și el, la rândul lui, un vector  $d$  – dimensional de forma:

$$m_x = \begin{bmatrix} m_1 \\ m_2 \\ \dots \\ m_d \end{bmatrix} \quad (5.189)$$

Componenta  $k$  a vectorului  $m_x$  este dată de:

$$m_k = E\{x_k\} = \int_{-\infty}^{+\infty} a_k f_x(a) da = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} a_k f_x(a) da_N \dots da_1 \quad (5.190)$$

După integrarea relației (5.190) pe toate componentele vectorului se obține:

$$m_k = \int_{-\infty}^{+\infty} a_k f_{x_k}(a_k) da_k \quad (5.191)$$

Dacă integrala (5.191) sau suma (în cazul exprimării relației (5.191) pentru o variabilă aleatoare discretă – relație similară (5.182)) nu converge către un număr real spunem că estimatorul statistic de ordinul întâi nu există.

### 3. Momente de ordin doi

În acest subcapitol am reunit generic, sub numele de momente de ordin doi (a se vedea și **Subcapitolul 5.5.2.1**):

- momentele de ordin doi ale unui vector aleator  $x$ ,
  - *matricea de corelație*,  $R_x$ ,
  - *matricea de covarianță*,  $C_x$ ;
- momentul centrat comun de ordin doi a două variabile aleatoare,
  - *covarianța între două variabile aleatoare*;
- momentele comune de ordin doi pentru doi vectori aleatori,  $x$  și  $y$ :
  - *matricea de cros-corelație*,  $R_{xy}$ ,
  - *matricea de cros-covarianță*,  $C_{xy}$ .

#### Matricea de corelație și de covarianță

**Matricea de corelație** reprezintă mulțimea completă de momente de ordin doi ale unui vector aleator,  $x$ , și ea este definită în modul următor:

$$R_x \stackrel{def}{=} E\{x x^{*T}\} \quad (5.192)$$

Pentru un vector aleator real obținem:

$$R_x = \begin{bmatrix} E\{(x_1)^2\} & E\{x_1 x_2\} & \cdots & E\{x_1 x_d\} \\ E\{x_2 x_1\} & E\{(x_2)^2\} & \cdots & E\{x_2 x_d\} \\ \vdots & \vdots & \ddots & \vdots \\ E\{x_d x_1\} & E\{x_d x_2\} & \cdots & E\{(x_d)^2\} \end{bmatrix} \quad (5.193)$$

Pentru un vector aleator complex matricea are forma:

$$R_x = \begin{bmatrix} E\{|x_1|^2\} & E\{x_1 x_2^*\} & \cdots & E\{x_1 x_d^*\} \\ E\{x_2 x_1^*\} & E\{|x_2|^2\} & \cdots & E\{x_2 x_d^*\} \\ \vdots & \vdots & \ddots & \vdots \\ E\{x_d x_1^*\} & E\{x_d x_2^*\} & \cdots & E\{|x_d|^2\} \end{bmatrix} \quad (5.194)$$

**Atenție!** În cazul determinării matricei de corelație (prin intermediul relațiilor (5.192) sau (5.194)) sau/și a matricei de covarianță (în modul în care va fi prezentată ulterior în cadrul acestui capitol) pentru cazul vectorilor aleatori complecși sau variabile aleatoare complexe este **obligatorie**, pe lângă utilizarea transpunerii, și utilizarea complex conjugatului în definirea tuturor acestor relații. În cazul neutilizării transpunerii complex conjugate, ci doar a

transpunerii simple în definierea matricei de corelație (relațiile (5.192) și (5.194)) vom obținem pentru orice vector aleator  $x$ :

$$E\{x x^T\} = [0] \quad (5.195)$$

Evident că rezultatul dat de relația (5.195) ( $R_x = [0]$ ), relație valabilă pentru orice vector aleator  $x$  complex, este unul care nu ne furnizează nici o informație relevantă despre vectorul  $x$ , această relație fiind astfel total inutilă în orice gen de analiză.

**Problema 5.25:** Demonstrați validitatea relației (5.195)<sup>49</sup>.

**Matricea de covarianță** este, la rândul ei, definită drept o mulțime de momente de ordinul doi centrate în jurul valorii medii. Această matrice este dată de relația:

$$C_x \stackrel{def}{=} E\{(x - m_x)(x - m_x)^{*T}\} \quad (5.196)$$

Forma matricială pentru  $C_x$  este similară cu cea prezentată în relațiile (5.193) și (5.194), cu mențiunea că, de această dată, elementele matricei sunt:

$$c_{kl} = \text{Cov}(x_k, x_l) = E\{(x_k - m_k)(x_l - m_l)^*\} \quad (5.197)$$

cu

$$c_{kk} = E\{|x_k - m_k|^2\} = \sigma_{x_k}^2 \quad (5.198)$$

Dacă matricea de covarianță se calculează pentru un vector de trăsături, atunci:

1. elementul  $c_{kl}$  este coeficientul de covarianță între trăsătura  $k$  și trăsătura  $l$  – notat în literatura de specialitate cu  $\text{Cov}(x_k, x_l)$ . Un coeficient de covarianță pozitiv ne arată că deviațiile de la valorile medii ale trăsăturilor  $x_k$  și  $x_l$  au, în medie, același semn. Dacă, creșterea (scăderea),

<sup>49</sup> Analiza vectorilor aleatori complecși,  $x = x_r + j \cdot x_i$  cu  $x_r$  și  $x_i$  vectori aleatori reali, (numai pentru cazul prelucrării digitale a semnalelor și pentru problemele de clasificare) poate fi redusă numai la acei vectori pentru care există următoarele relații de simetrie [Therrien, 1992]:

$$E\{x_r x_r^T\} = E\{x_i x_i^T\} \quad (a)$$

$$E\{x_i x_r^T\} = -E\{x_r x_i^T\} \quad (b)$$

Aceste două relații ne furnizează informația că părțile reale și imaginare ale vectorului satisfac condițiile:

$$E\{x_{rk} x_{rl}\} = E\{x_{ik} x_{il}\} \quad (b)$$

și

$$E\{x_{ik} x_{rl}\} = -E\{x_{rk} x_{il}\} \quad (c)$$

în medie, a valorii unei trăsături față de media ei va determina, în medie, și creșterea (scăderea) valorii celei de a doua trăsături față de media proprie, atunci spunem că trăsăturile sunt **corelate pozitiv**. O valoare negativă a acestui coeficient ne indică că, în medie, deviațiile acestor două trăsături au sensuri opuse față de mediile lor și, în această situație trăsăturile spunem că sunt **corelate negativ**. Dacă acest coeficient este zero spunem că variabilele aleatoare  $x_k$  și  $x_l$  sunt **necorelate**. În situația în care  $x_k$  și  $x_l$  sunt **statistic independente** atunci  $c_{kl} = 0$ .

2.  **$c_{kk}$  este varianța trăsăturii  $k$**  – acest parametru furnizează o informație despre „întinderea”, dispersia spațială, a valorilor trăsăturii  $k$ .

În continuare exemplificăm forma matricii de covarianță pentru cazurile unor distribuții unidimensionale (1D), bidimensionale (2D) și tridimensionale (3D). Pentru astfel de vectori:

$$x_{(1D)} = [x_1] = x_1 \quad x_{(2D)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x_{(3D)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad (5.199)$$

matricile de covarianță (pentru cele trei situații prezentate anterior) iau următoarele forme:

$$\begin{aligned} C_{x(1D)} &= [E\{(x_1 - m_{x_1})(x_1 - m_{x_1})^*\}] = \\ &= [E\{|x_1 - m_{x_1}|^2\}] = [\sigma_{x_1}^2] = \sigma_{x_1}^2 \end{aligned} \quad (5.200)$$

$$\begin{aligned} C_{x(2D)} &= \begin{bmatrix} E\{(x_1 - m_{x_1})(x_1 - m_{x_1})^*\} & E\{(x_1 - m_{x_1})(x_2 - m_{x_2})^*\} \\ E\{(x_2 - m_{x_2})(x_1 - m_{x_1})^*\} & E\{(x_2 - m_{x_2})(x_2 - m_{x_2})^*\} \end{bmatrix} = \\ &= \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix} = \begin{bmatrix} \sigma_{x_1}^2 & Cov(x_1, x_2) \\ Cov(x_2, x_1) & \sigma_{x_2}^2 \end{bmatrix} \end{aligned} \quad (5.201)$$

$$\begin{aligned} C_{x(2D)} &= \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix} = \\ &= \begin{bmatrix} \sigma_{x_1}^2 & Cov(x_1, x_2) & Cov(x_1, x_3) \\ Cov(x_2, x_1) & \sigma_{x_2}^2 & Cov(x_2, x_3) \\ Cov(x_3, x_1) & Cov(x_3, x_2) & \sigma_{x_3}^2 \end{bmatrix} \end{aligned} \quad (5.202)$$

Generalizând relațiile (5.200), (5.201) și (5.202) constatăm că pe diagonala principală vom avea întotdeauna varianța trăsăturilor în timp ce restul termenilor sunt dați de coeficienții de covarianță – având semnificația

varianței relative, a varianței unei variabile (de exemplu,  $x$  sau  $x_1$ ) relativă la varianța unei alte variabile aleatoare (de exemplu,  $y$  sau  $x_3$ ). Din acest motiv numele corect al acestei matrici este: **matricea de varianță-covarianță** (a unui anumit vector aleator). Cu toate acestea numele care s-a impus este acela de **matrice de covarianță** (potențial datorită numărului superior de termeni de covarianță, față de cei de varianță, începând de la matrici cu dimensionalitate minimă de  $3 \times 3$ ).

**Problema 5.26:** Știind că media unei variabile aleatoare reale este  $m_x = E\{x\}$  în timp ce varianța aceleiași variabile este dată de  $\sigma_x^2 = E\{(x - m_x)^2\}$  determinați  $E\{x^2\}$  în funcție numai de  $m_x$  și  $\sigma_x^2$ .

Din relațiile (5.192) și (5.123) se poate observa că ambele matrici, atât cea de corelație cât și cea de covarianță sunt **matrici Hermetiene simetrice**, adică au proprietatea:

$$R_x = R_x^{*T} \quad (5.203)$$

și

$$C_x = C_x^{*T} \quad (5.204)$$

**Problema 5.27:** Demonstrați validitatea relațiilor (5.203) și (5.204).

Mai mult decât atât, ambele matrici sunt și **pozitiv semidefinite**, deci satisfac relațiile:

$$a^{*T} R_x a \geq 0 \quad (5.205)$$

și

$$a^{*T} C_x a \geq 0 \quad (5.206)$$

pentru orice vector complex  $a$ . Relațiile (5.205) și (5.206) se pot demonstra foarte ușor, ținând cont că<sup>50</sup>:

$$a^{*T} R_x a = a^{*T} E\{xx^{*T}\} a = E\{|x^{*T}a|^2\} \quad (5.207)$$

este întotdeauna mai mare sau egală cu zero.

Matricea de covarianță și cea de corelație sunt legate între ele printr-o relație matematică. Observând că  $E$  este un operator liniar iar  $m_x$  este un vector având elemente constante putem scrie:

$$\begin{aligned} E\{(x - m_x)(x - m_x)^{*T}\} &= E\{xx^{*T} - xm_x^{*T} - m_x x^{*T} + m_x m_x^{*T}\} = \\ &= E\{xx^{*T}\} - E\{x\}m_x^{*T} - m_x E\{x^{*T}\} + m_x m_x^{*T} \quad (5.208) \\ &= E\{xx^{*T}\} - m_x m_x^{*T} \end{aligned}$$

Rezultând **o relație foarte importantă**:

<sup>50</sup> Pentru demonstrarea relațiilor (5.205) și (5.206) țineți cont și de relația:  $(A B^T)^T = B A^T$

$$R_x = C_x + m_x m_x^{*T} \quad (5.209)$$

### Covarianța între două variabile aleatoare

Așa după cum am văzut anterior, un element al matricei de covarianță, care nu este poziționat pe diagonala principală, este caracterizat de relația:

$$c_{kl} = E\{ (x_k - m_k)(x_l - m_l)^* \} \quad (5.210)$$

Această valoare **măsoară tendința trăsăturilor**  $x_k$  și  $x_l$  **de a varia împreună**.

După cum se va prezenta ulterior, dacă  $a^1, a^2, \dots, a^K$  sunt  $K$  realizări particulare ale vectorului aleator de trăsături  $x$ , atunci relația de estimare a covarianței între trăsăturile  $k$  și  $l$  ale vectorului  $x$  (caracterizate de variabilele aleatoare  $x_k$  și  $x_l$ ) se calculează cu relația (vezi **Observația 5.27**):

$$\hat{c}_{kl} = \frac{1}{K} \left[ (a_k^1 - \hat{m}_k)(a_l^1 - \hat{m}_l)^* + (a_k^2 - \hat{m}_k)(a_l^2 - \hat{m}_l)^* + \dots + (a_k^K - \hat{m}_k)(a_l^K - \hat{m}_l)^* \right] \quad (5.211)$$

În relația anterioară  $\hat{m}_z$ , cu  $z$  egal cu  $k$  sau  $l$ , este estimatorul mediei trăsăturii  $k$  sau  $l$ . Dacă realizările particulare ale vectorului aleator  $x$  sunt vectori reali, complex conjugatul nu-și mai are rolul în relația anterioară.

**Observația 5.40:** Dacă, în contextul de mai sus, dorim să determinăm un coeficient de covarianță oarecare,  $c_{kl}$ , în alt mod decât cel prezentat mai sus, atunci, putem estima mai întâi matricea de covarianță a setului de date după care, ulterior preluăm din aceasta spre analiză numai valoarea elementului care ne interesează.

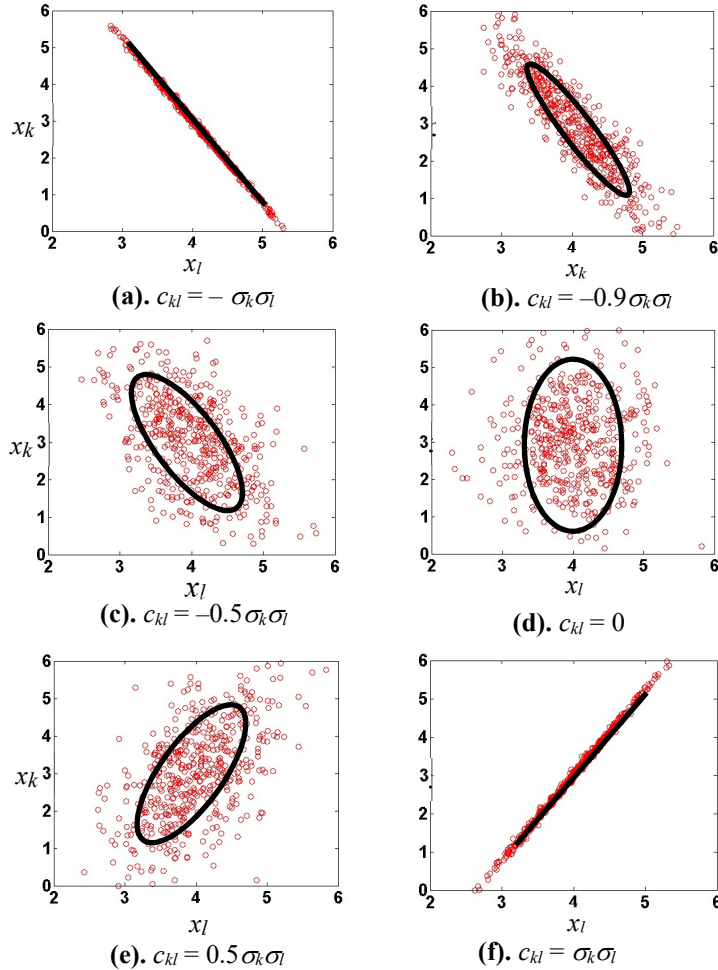
*Coeficientul de covarianță între două variabile aleatoare oarecare* – fie aceste componente de ordin  $k$  și  $l$  ale unui vector aleator de trăsături  $x$  – are, așa cum am mai amintit, următoarele **proprietăți fundamentale**:

- dacă trăsătura  $k$  și trăsătura  $l$  au tendința, în medie, de a varia în același sens, atunci  $c_{kl} > 0$ ;
- dacă, în medie, valoarea trăsăturii  $k$  crește iar cea a trăsăturii  $l$  scade (sau invers) față de mediile lor corespunzătoare, atunci  $c_{kl} < 0$ ;
- în cazul în care trăsăturile sunt independente avem  $c_{kl} = 0$ ;
- $|c_{kl}| \leq \hat{\sigma}_k \cdot \hat{\sigma}_l$ , unde  $\hat{\sigma}_k$  și  $\hat{\sigma}_l$  sunt deviațiile standard estimate pentru cele două trăsături,  $k$  și, respectiv,  $l$ ;
- $c_{kk} = \hat{\sigma}_k^2$  este varianța estimată a trăsăturii  $k$ .

Proprietățile enumerate anterior pentru coeficientul de covarianță sunt adevărate numai dacă anumite condiții sunt satisfăcute. Astfel, în cazul



trăsăturilor necorelate,  $c_{kl}$  este egal cu zero numai dacă există un număr infinit de realizări particulare ale variabilei aleatoare  $x$ ; în caz contrar (număr mare dar finit de exemplare) această valoare nu va fi zero ci va fi o valoare foarte mică, apropiată de zero. Punerea în evidență a influenței mărimii setului de date asupra diferiților parametri estimați ai matricei de covarianță va fi evidențiată în **Capitolul 7**.

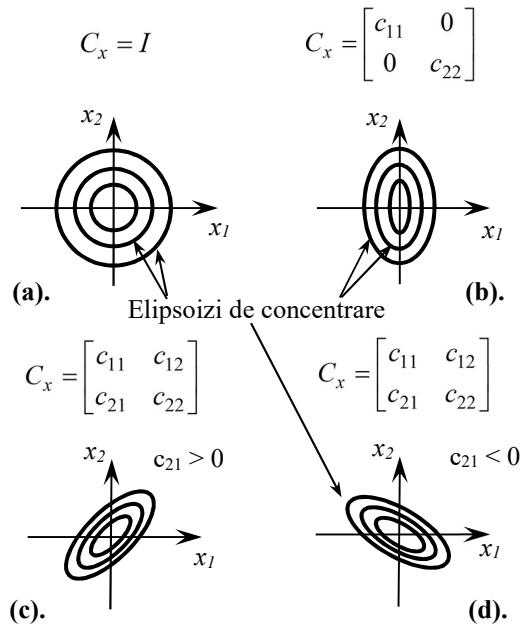


**Figura 5.23.** Relația ce există între diferite valori ale coeficientului de covarianță și distribuția reală a elementelor în spațiul trăsăturilor

Pentru o problemă reală, caracterizată de un număr suficient de mare de eșantioane, valoarea lui  $c_{kl}$ , pentru situația trăsăturilor dependente, va fi cu mult mai mică decât valoarea maximă pe care o poate lua acest coeficient, respectiv valoare egală cu produsul deviațiilor standard ale trăsăturilor  $k$  și  $l$ .

Covarianța a două variabile aleatoare,  $x_k$  și  $x_l$ , ce descriu două trăsături ale unei populații poate lua valori în intervalul  $[-\sigma_k \cdot \sigma_l, +\sigma_k \cdot \sigma_l]$  și ea măsoară dependența/independența între cele două variabile.

Correspondența între valoarea coeficientului de covarianță între două trăsături și o posibilă distribuție a elementelor unei clase funcție de aceste două trăsături este dată în **Figura 5.23**. Elipsele prezentate în **Figura 5.23** sunt contururile funcției densitate de probabilitate (intersecția dintre funcția densitate de probabilitate și diferite planuri paralele cu planul trăsăturilor, vezi **Figura 5.7**) proiectate în planul trăsăturilor ( $x_k$  versus  $x_l$ ). Aceste elipse poartă numele de *elipse de concentrare*. În **Figura 5.24** se prezintă legăturile care există între diferite reprezentări ale elipselor de concentrare și posibilele matrici de covarianță generatoare ale acestor reprezentări.

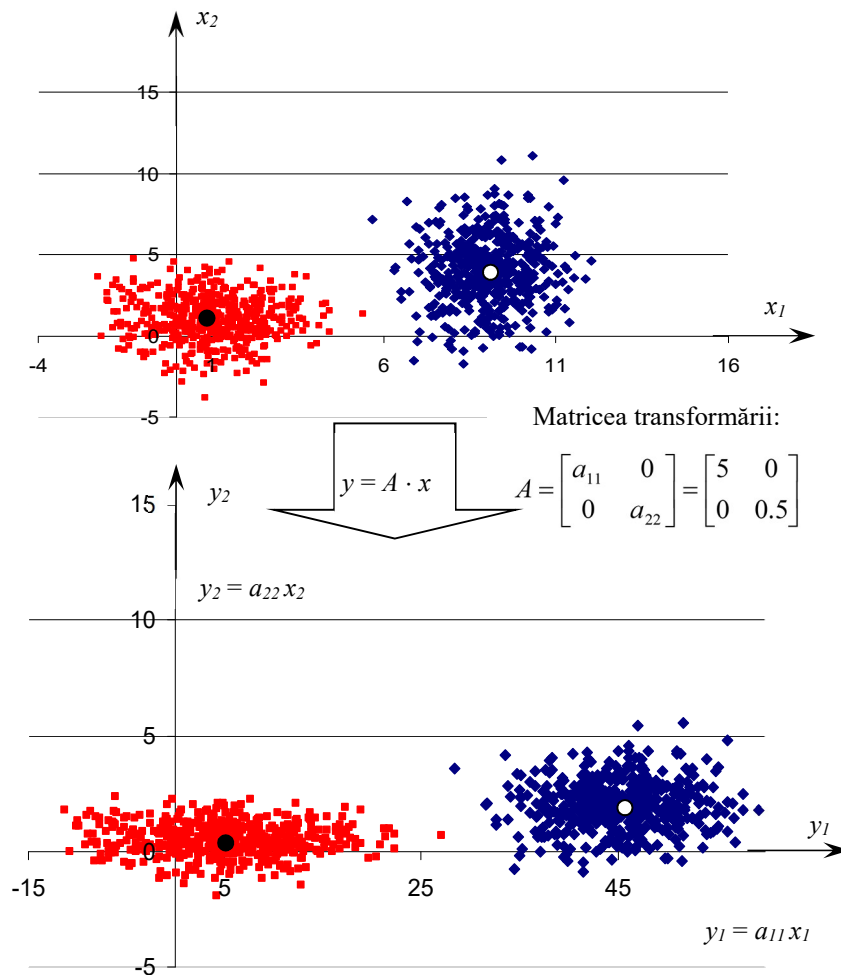


**Figura 5.24.** Reprezentarea diferitelor distribuții Gauss-iene bidimensionale prin intermediul proiecției cotelor de egală probabilitate – coeficienții  $c_{kl}$  sunt dați de relațiile (5.197) și (5.198)

**Aplicația 5.2:** Utilizând programul din directorul **TrasareDensitateGauss** asociat acestui capitol:

1. Introduceți diferite valori pentru coeficienții matricii  $C_x$  astfel încât să obțineți în mod calitativ toate cazurile prezentate în **Figura 5.24**.

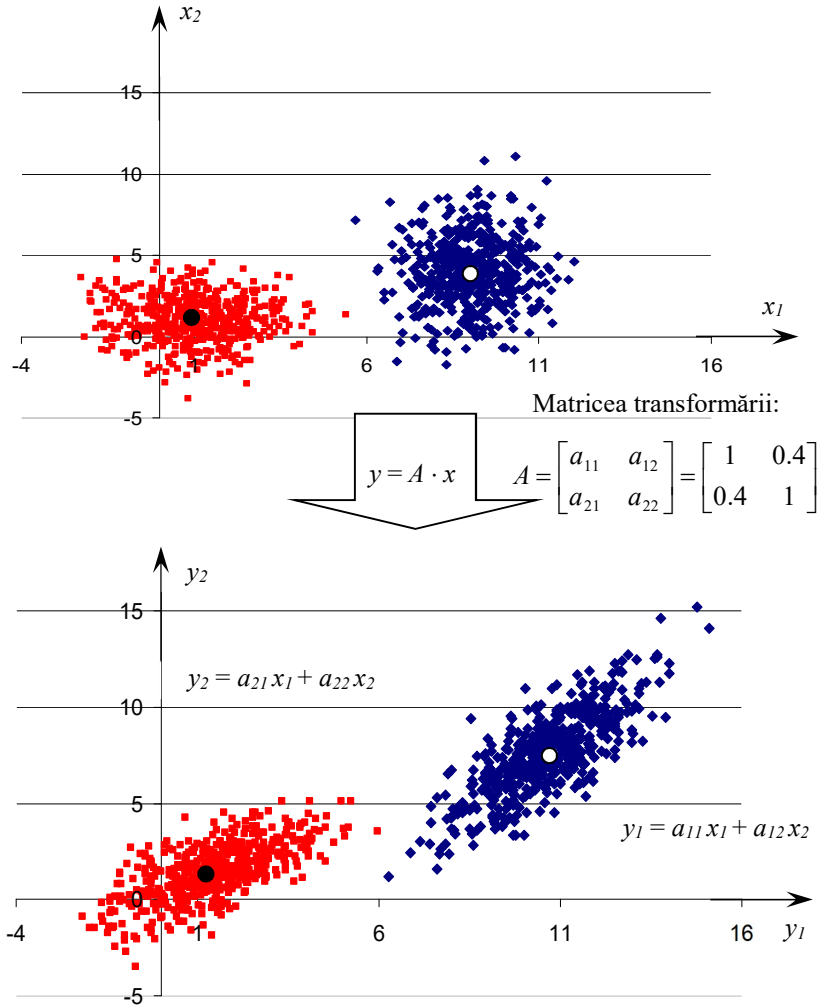
2. Introduceți diferite numere întregi în *EditBox*-ul „*Levels*” pentru a obține diferite contururi de echiprobabilitate ale funcției densitate de probabilitate *Gauss*-iană. Modificați numărul de puncte pe care faceți analiza și rotiți reprezentările pentru îmbunătățirea vizibilității. Pentru a observa simultan atât funcția densitate de probabilitate cât și repartiția punctelor în plan modificați transparența funcției densitate de probabilitate gauss-iene reprezentată grafic prin intermediul *EditBox*-ului „*Transparency*”.



**Figura 5.25.** Efectul scalării asupra distribuțiilor celor două clase

**Întrebări:**

1. Cum trebuie să fie coeficienții  $c_{11}$  și  $c_{22}$  astfel încât o elipsă de concentrare a distribuției *Gauss*-iene să aibă axa mare paralelă cu axa  $x_1$  a sistemului de coordonate (vezi comparativ și **Figura 5.24 (b)**)?
2. Coeficienții  $c_{21}$  și  $c_{12}$  pot lua maxim valoarea  $\hat{\sigma}_1 \cdot \hat{\sigma}_2$ . Atribuiți această valoare, cu semnul „-” sau „+”, acestor coeficienți. Explicați rezultatele obținute și vizualizate.



**Figura 5.26.** Un posibil efect al unei transformări generale asupra distribuțiilor a două clase

**Observația 5.41:** Diferența conceptuală care există între **Figura 5.24** și

**Figura 5.23** este dată în principal de faptul că în ultimul set de figuri este tratat cazul unui vector de trăsături bidimensional și al corelației ce există între cele două trăsături ale lui, în timp ce în **Figura 5.23** este prezentat un caz mai general și posibilele situații existente între oricare două trăsături (de exemplu între trăsăturile  $k$  și  $l$ ) ale aceluiași vector care de această dată poate să fie unul  $d$  dimensional.

Există o serie largă de mecanisme care pot determina corelarea trăsăturilor unui vector aleator. De exemplu, aplicarea unei **transformări liniare** asupra setului de date poate avea un asemenea efect (pentru o tratare mai în detaliu a transformărilor liniare vezi **Subcapitolele 6.1 și 6.2**).

*Operația de scalare* este un exemplu particular al **transformărilor liniare**. Din punct de vedere geometric operația de scalare, „dilată” sau „contractă” pe una sau pe mai multe direcții pe care se aplică aceasta (care corespund cu una sau mai multe axe ale sistemului de coordonate) forma aglomerărilor de realizări particulare ale unei variabilei aleatoare. Astfel, de exemplu, clusteri de vectori care inițial erau circulari devin de formă elipsoidală, cu axele principale ale elipsoizilor de concentrare ce caracterizează aceste distribuții orientate paralel cu axele sistemului de coordonate.

O transformare liniară mai generală introduce în plus și o rotire în planul trăsăturilor. Astfel, de exemplu, clusterii care inițial erau circulari devin acum elipsoidali cu axele principale ce fac un anumit unghi față de axele sistemului de coordonate, **Figura 5.26**. În acest mod se introduce o covarianță între componentele vectorului de trăsături.

**Exemplu 5.13:** Dacă considerăm situația achiziționării simultane, pe  $N$  canale, a semnalului EEG (în multe din situațiile practice  $N$  poate lua și valori mai mari de 20), atunci acestea pot fi grupate sub forma unui **vector aleator**,  $x = [x_1, \dots, x_N]^T \in R^N$ , ce reprezintă o secvență de  $N$  variabile aleatoare. Deoarece conductorii pe care se transmite semnalul EEG sunt foarte apropiați, fiecare canal induce în toate celelalte un anumit semnal perturbator prin intermediul cuplajului capacitiv (prin intermediul câmpului electric) și prin intermediul cuplajului inductiv (prin intermediul câmpului magnetic)<sup>51</sup>. Astfel, la intrarea sistemului de achiziție a unui canal vom avea informația de la nivelul electrodului ce corespunde aceluși canal plus o combinație de semnale (considerate aici perturbatoare) de la celelalte 19 canale:

$$y_i = a_{i1} x_1 + a_{i2} x_2 + \dots + a_{iN} x_N \quad (5.212)$$

<sup>51</sup> Acest fenomen poartă denumirea de **diafonie**.

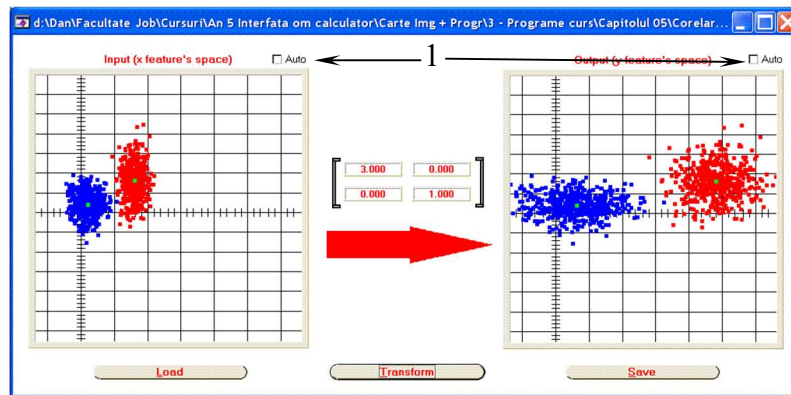
Global pentru toate cele  $N$  canale de achiziție putem scrie:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{NN} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \quad (5.213)$$

Se observă că *transformarea liniară* prezentată în (5.213) este obținută prin simpla transmisie a informațiilor de la electrozii de achiziție la intrările sistemului de conversie analog-numerică. Această simplă transmisie a datelor determină o corelare mai mică sau mai mare a setului de date de intrare – achiziționat de la senzorii poziționați pe scapul subiectului.

**Aplicația 5.3:** Utilizând programul din directorul *Corelare seturi de date 2D* asociat acestui capitol precum și setul de date „Date 2 clase necorelate B.txt”, parcurgeți și răspundeți cerințelor de la următoarele puncte:

1. Aplicați transformările prezentate în **Figura 5.25** și **Figura 5.26** și obțineți transformări similare ale setului de date, conforme cu figurile anterioare.
2. Determinați o transformată  $A$  (vezi **Figura 5.25** și **Figura 5.26**) astfel încât să obțineți o corelare „negativă” a setului de date – caracterizată de o distribuție a elementelor conformă cu **Figura 5.23 (b)** sau (c).



**Figura 5.27.** Interfața grafică a programului

3. Intuitiv, ținând constante valorile elementelor  $a_{11}$  și  $a_{22}$  din matricea  $A$  și considerând egale valorile elementelor  $a_{12}$  și  $a_{21}$ , în care din cele două

situații valorile elementului  $a_{12}$  va fi mai mic dacă se dorește obținerea unor distribuții similare cu cele din **Figura 5.23(b)** sau **(c)**?

Programul de mai sus este capabil să reprezinte atât spațiul de intrare cât și cel de ieșire, păstrând aceleași unități de măsură și intervale pentru ambele axe astfel încât să se poată face o comparație corectă a distribuției elementelor în cele două spații. În momentul în care se dorește o scalare a seturilor de date pe cele două clase se ține cont și de valorile maxime ale elementelor se utilizează check-boxurile 1, vezi **Figura 5.27**. Încărcarea setului de date de intrare, salvarea rezultatelor și aplicarea transformării liniare,  $A$ , se inițiază din butoanele poziționate în partea de jos a aplicației, **Figura 5.27**.

### **Matricea de cros-corelație și cros-covarianță**

Când o cantitate  $g(\cdot)$  depinde de doi vectori aleatori  $x$  și  $y$  (cum este cazul, de exemplu, al momentelor comune pentru doi vectori aleatori), media acesteia poate fi scrisă:

$$E\{g(x, y)\} = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(a, b) f_{xy}(a, b) da db \quad (5.214)$$

Momentele statistice comune cele mai folosite pentru doi vectori aleatori oarecare sunt **matricea de cros-corelație**:

$$R_{xy} = E\{x y^{*T}\} \quad (5.215)$$

și **matricea de cros-covarianță**:

$$C_{xy} = E\{(x - m_x)(y - m_y)^{*T}\} \quad (5.216)$$

În cazul cel mai general nici una dintre matricile  $R_{xy}$  sau  $C_{xy}$  nu este pătratică. Acestea pot fi însă pătratice numai în situația în care  $x$  și  $y$  au același număr de elemente. Printr-o demonstrație similară cu cea dată în momentul obținerii relației (5.209) se poate arăta că:

$$R_{xy} = C_{xy} + m_x m_y^{*T} \quad (5.217)$$

Deși ambele matrici  $R_{xy}$ ,  $C_{xy}$  nu sunt, în general, simetrice se poate arăta că:

$$R_{xy} = R_{yx}^{*T} \quad (5.218)$$

și

$$C_{xy} = C_{yx}^{*T} \quad (5.219)$$

**Problema 5.28:** Demonstrați validitatea relațiilor (5.218) și (5.219).

### **Proprietăți ale momentelor statistice de ordin doi**

Pentru anumite valori particulare ale momentelor statistice de ordin doi acestea caracterizează, în mod specific, vectorii aleatori asociați. Astfel:

(a). Vectorii aleatori  $x$  și  $y$  sunt **necorelați** dacă:

$$R_{xy} = E\{x y^{*T}\} = E\{x\} E\{y^{*T}\} = m_x m_y^{*T} \quad (5.220)$$

Datorită relației (5.217), putem afirma, în mod echivalent, că vectorii aleatori sunt **necorelați** dacă:

$$C_{xy} = E\{(x - m_x)(y - m_y)^{*T}\} = [0] \quad (5.221)$$

deci, dacă **matricea de cros-covarianță este zero** – aceasta este, de altfel, și cea mai uzuală definiție.

(b). Doi vectori sunt **ortogonali**, dacă:

$$R_{xy} = E\{x y^{*T}\} = [0] \quad (5.222)$$

deci, dacă **matricea de cros-corelație este zero**.

Din relația (5.217) putem vedea că în momentul în care media ambilor vectori este zero atunci ambele matrici, de cros-corelație și, respectiv, de cros-covarianță sunt similare; numai în această situație vectorii aleatori ortogonali sunt și necorelați și invers.

(c). Din relația (5.220) este ușor de demonstrat că **atunci când vectorii aleatori  $x$  și  $y$  sunt independenți<sup>52</sup> atunci sunt și necorelați**. Inversa acestei afirmații este în general falsă. Putem spune că această condiție de independență este o condiție mai „dură”, mai restrictivă decât cea de necorelare deoarece independența implică automat și necorelarea.

**Problema 5.29:** Încercați să demonstrați inversa afirmației precedente, atunci când cei doi vectori sunt caracterizați de o distribuție *Gauss-iană* de probabilitate.

**Problema 5.30:** Demonstrați că atunci când vectorii aleatori  $x$  și  $y$  sunt independenți atunci sunt și necorelați.

Aceste concepte de ortogonalitate și necorelație sunt definite în mod similar și pentru componentele  $x_k$  și  $x_l$  ale unui vector aleator  $x$ . Astfel:

<sup>52</sup> Doi vectori aleatori  $x$  și  $y$  sunt statistic independenți dacă:  $f_{xy}(x,y) = f_x(x) \cdot f_y(y)$



(d). **Componentele  $k$  și  $l$**  ale vectorului aleator  $x$  vor fi **ortogonale** dacă:

$$E\{x_k x_l^{*T}\} = 0; \quad k \neq l \quad (5.223)$$

Ca o consecință directă, matricea de corelație,  $R_x$ , a vectorului aleator  $x$  este o matrice diagonală atunci când toate componentele vectorului sunt ortogonale două câte două.

(e). Aceleași **componente vor fi necorelate** dacă:

$$E\{(x_k - m_k)(x_l - m_l)^{*T}\} = 0; \quad k \neq l \quad (5.224)$$

În mod echivalent putem spune că **matricea de covarianță,  $C_x$ , a vectorului aleator  $x$  este o matrice diagonală atunci când componentele vectorului aleator sunt necorelate** (două câte două).

Atenție! Deosebirea majoră între proprietățile (a), (b) și proprietățile (d) și (e) este dată de faptul că primele două proprietăți se referă la doi vectori aleatori diferiți în timp ce proprietățile (d) și (e) se referă la două componente (variabile aleatoare) ale aceluiași vector aleator.

**Observația 5.42:** Este util de știut faptul că, dându-se un vector aleator  $x$  cu orice tip de distribuție, se poate găsi o transformare liniară din spațiul inițial în alt spațiu astfel încât componentele vectorului, în noul spațiu, să fie ortogonale și/sau necorelate. Astfel de tehnici vor fi prezentate în capitolele următoare. Folosindu-ne de o astfel de transformare liniară putem depăși cu succes neajunsul prezentat în **Subcapitolul 4.3.2**, respectiv prezența corelației între diferitele componente ale vectorului de trăsături; în acest context, în cazul seturilor de date caracterizate de trăsături puternic corelate inclusiv clasificatorul de tip minimă distanță poate furniza performanțe superioare cu condiția, bineînțeles, a unei pre-procesări prealabile constând în decorelarea datelor.

### 5.5.3. Estimarea parametrilor statistici din setul de date

În acest subcapitol se prezintă modul în care se pot estima parametrii statistici (medie, matrice de covarianță etc.) a unei populații sau proces atunci când funcția densitate de probabilitate este necunoscută, dar se cunosc, în schimb, o serie de realizări particulare ale vectorului aleator  $x$ :  $a^1, a^2, \dots, a^K$ .

Dacă ar fi să determinăm media, matricea de corelație sau cea de covarianță a unui vector aleator  $x$  atunci, conform modului cum au fost definite aceste momente (vezi **Subcapitolul 5.5.2**, subpunctul 1 „**Momentele în analiza statistică. Funcții caracteristice**”), avem nevoie să cunoaștem, în

prealabil, funcția densitate de probabilitate a lui  $x$ . Reluând mai jos relațiile prin care sunt definite aceste cantități, avem că:

$$m_x = E\{x\} \stackrel{def.}{=} \int_{-\infty}^{+\infty} a f_x(a) da \quad (5.225)$$

$$R_x = E\{xx^{*T}\} = \int_{-\infty}^{+\infty} aa^{*T} f_x(a) da \quad (5.226)$$

$$C_x = E\{(x - m_x)(x - m_x)^{*T}\} = \int_{-\infty}^{+\infty} (a - m_x)(a - m_x)^{*T} f_x(a) da \quad (5.227)$$

În continuare prezentăm și demonstrăm modul în care se pot estima, din setul de date, diferiți parametri statistici fără a cunoaște efectiv funcția densitate de probabilitate. Această prezentare va fi făcută pentru un vector aleator real monodimensional – deci, vom lua în considerare cazul unei variabile aleatoare simple. Generalizarea acestei demonstrații (de la situația monodimensională la o alta, în mai multe dimensiuni) este una directă.

### 1. Estimarea momentelor statistice pentru un vector aleator

Fie  $g(x)$  o funcție ce are ca argument vectorul aleator  $x$ . Dacă ambele mărimi  $f_x(a)$  și  $g(a)$  nu variază substanțial pe toate intervalele definite, atunci, folosind relația (5.166), putem realiza aproximarea:

$$E\{g(x)\} = \int_{-\infty}^{+\infty} g(a) f_x(a) da \cong \sum_{\text{toti } a^0} g(a^0) f_x(a^0) \Delta a \cong \frac{1}{K} \sum_{\text{toti } a^0} g(a^0) K_{a^0} \quad (5.228)$$

unde  $K$  reprezintă numărul total de valori ale lui  $x$  dintr-un eșantion iar  $K_{a^0}$  este numărul de valori din eșantion aflate în intervalul  $(a^0, a^0 + \Delta a]$ ; deci, pe un singur interval de lungime  $\Delta a$  care pornește din punctul  $a^0$  avem o mulțime formată din  $K_{a^0}$  elemente, respectiv:

$$\left\{ a^i, a^{i+1}, \dots, a^{i+K_{a^0}-1} \right\} \quad (5.229)$$

Pentru toți vectorii  $a^i$  aparținând aceluiași interval, în ipoteza  $g(a^i) \approx g(a^0)$  (ipoteză necesară aproximării integralei din relația (5.228)), cantitatea  $g(a^0)K_{a^0}$  poate fi scrisă sub forma:

$$K_{a^0} \cdot g(a^0) = \underbrace{g(a^0) + g(a^0) + \dots + g(a^0)}_{\text{de } K_{a^0} \text{ ori}} = \sum_{j=i}^{i+K_{a^0}-1} g(a^j) \quad (5.230)$$

Relația (5.228) poate fi rescrisă prin sumarea tuturor cantităților de tipul (5.230). Dacă această procedură este urmată pentru toate intervalele, relația (5.228) devine:

$$E\{g(x)\} \cong \frac{1}{K} \sum_{k=1}^K g(a^k) \quad (5.231)$$

În acest mod am obținut un *estimat empiric* (valoare aproximată) pentru media vectorului/matricii aleator(oare)  $g(x)$  a lui  $x$ ; în plus, se poate observa, din relația de mai sus, că **estimatorul** este dat chiar de **media aritmetică a eșantionului**, calculată pentru  $g(x)$ .

Particularizând definiția funcției  $g(x)$ , respectiv,  $g(x) = x$ ,  $g(x) = x x^{*T}$  și  $g(x) = (x - m_x)(x - m_x)^{*T}$ , vom obține următoarele estimări pentru media vectorului aleator  $x$ , matricea de corelație și, respectiv, matricea de covarianță a lui  $x$ :

$$\hat{m}_x = E\{g(x)\} = E\{x\} = \frac{1}{K} \sum_{k=1}^K a^k \quad (5.232)$$

$$\hat{R}_x = E\{g(x)\} = E\{x \cdot x^{*T}\} = \frac{1}{K} \sum_{k=1}^K a^k (a^k)^{*T} \quad (5.233)$$

$$\hat{C}_x = E\{g(x)\} = E\{(x - m_x) \cdot (x - m_x)^{*T}\} = \frac{1}{K} \sum_{k=1}^K (a^k - m_x)(a^k - m_x)^{*T} \quad (5.234)$$

## 2. Estimarea momentelor statistice comune pentru doi vectori aleatori

Un rezultat similar se obține și în cazul estimării cantităților  $g(x,y)$ , dependente de doi vectori aleatori,  $x$  și  $y$  – vectori ce sunt caracterizați prin realizările lor particulare:  $a^1, a^2, \dots, a^k$  și, respectiv,  $b^1, b^2, \dots, b^k$ . În această situație un estimat pentru parametrul *media* vectorului/matricii aleator(oare)  $g(x,y)$  este dat de relația:

$$E\{g(x, y)\} \cong \frac{1}{K} \sum_{k=1}^K g(a^k, b^k) \quad (5.235)$$

unde estimatorul este, de asemenea, *media aritmetică a eșantionului*, calculată însă, de această dată, pentru  $g(x, y)$ .

Și de această dată, particularizând definiția funcției  $g(x, y)$ , respectiv,  $g(x, y) = xy^{*T}$  și  $g(x, y) = (x - m_x)(y - m_y)^{*T}$ , obținem următoarele estimări pentru matricile de cros-corelație, respectiv, cros-covarianță ale vectorilor aleatori  $x$  și  $y$ :

$$\hat{R}_{xy} = E\{g(x, y)\} = E\{x \cdot y^{*T}\} = \frac{1}{K} \sum_{k=1}^K a^k (b^k)^{*T} \quad (5.236)$$

$$\hat{C}_{xy} = E\{g(x, y)\} = E\{(x - m_x) \cdot (y - m_y)^{*T}\} = \frac{1}{K} \sum_{k=1}^K (a^k - m_x) (b^k - m_y)^{*T} \quad (5.237)$$

### 3. Estimarea momentelor statistice din matricea setului de date

Pentru estimarea, de exemplu, a matricelor de corelație și cros-corelație a unui vector aleator  $x$ , respectiv, a doi vectori aleatori,  $x$  și  $y$ , se pot folosi, așa cum am arătat mai sus, în mod direct, relațiile (5.231) și, corespunzător, (5.235). La același rezultat se poate ajunge, însă, și printr-o altă metodă echivalentă, metodă pe care o prezentăm în cele ce urmează.

Astfel, într-un prim pas se definește matricea vectorilor de date sub forma:

$$X = \begin{bmatrix} - (a^1)^{*T} & - \\ - (a^2)^{*T} & - \\ \vdots & \\ - (a^K)^{*T} & - \end{bmatrix} \quad (5.238)$$

Matricea  $X$  mai poartă numele de *matricea eșantionului* sau *matricea setului de date*.

Folosind următorul produs:

$$X^{*T} X = \sum_{k=1}^K a^k (a^k)^{*T} \quad (5.239)$$

putem estima, de exemplu, matricea de corelație, *în mod direct*, prin relația:

$$\hat{R}_x = \frac{1}{K} X^{*T} X \quad (5.240)$$

Formulele de calcul pentru estimarea matricii de corelație, date de relațiile (5.233) și (5.240) sunt, deci, echivalente, cu mențiunea că ultima dintre acestea este, însă, una mult mai directă.

*Estimarea matricii de covarianță* se poate face într-un mod similar celui prezentat mai sus. Singura deosebire este, în acest caz, dată de maniera în care definim, pentru început, matricea de calcul. Astfel, în locul matricii setului de date,  $X$ , construim matricea  $X_0$ , dată de elementele eșantionului din care s-a extras, mai întâi, vectorul medie  $m_x$ :

$$a_0^i = a^i - m_x, \quad \text{pentru } i = \overline{1, K} \quad (5.241)$$

$$X_0 = \begin{bmatrix} - & (a_0^1)^{*T} & - \\ - & (a_0^2)^{*T} & - \\ & \vdots & \\ - & (a_0^K)^{*T} & - \end{bmatrix} \quad (5.242)$$

În continuare, pentru determinarea matricii de covarianță aplicăm relația:

$$\hat{C}_x = \frac{1}{K} X_0^{*T} X_0 \quad (5.243)$$

Modalitatea de estimare a matricii de covarianță, dată de relația (5.243), implică calcularea matricii  $X_0$  prin eliminarea din fiecare vector de trăsături în parte a valorii medii (conform relației (5.241)) și apoi se va compune matricea  $X_0$  din acești vectori. Această abordare este una ușor de abordat pentru un operator uman. Un sistem informatic (embedded – DSP, microcontroler etc.) va lucra în mod direct cu matrici, pornind de la matricea setului de date, relația (5.238), și prin utilizarea unor calcule matriciale este de dorit a obține direct matricea  $X_0$ . O astfel de abordare se prezintă mai departe:

$$X_0 = X - \frac{1}{K} \cdot [1] \cdot X \quad (5.244)$$

În relația (5.244) matricea  $[1]$  este o matrice pătratică, având o dimensionalitate  $K \times K$  ( $K$  este numărul de vectori de trăsături utilizați în estimarea matricii de covarianță), cu toate elementele egale cu 1.

**Problema 5.31:** Demonstrați corectitudinea relației (5.244).

**Problema 5.32:** Pentru următoarele realizări particulare ale unui vector aleator  $x$  bidimensional,  $a^1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ ,  $a^2 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ ,  $a^3 = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$ ,  $a^4 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$  estimați:

1. Media vectorului aleator, ce caracterizează clasa generată de realizările particulare anterior prezentate.
2. Matricea de covarianță a clasei (estimarea matricei de covarianță se va realiza prin două metode diferite).

**Rezolvare:**

1. Pentru estimarea vectorului mediu al clasei ne folosim de relația (5.232) particularizată după cum urmează:

$$\hat{m}_x = \frac{1}{4} \sum_{i=1}^4 a^i = \frac{1}{4} \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 3 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 \\ 2 \end{bmatrix} \right\} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}.$$

2. Pentru estimarea matricei de covarianță vom extrage, într-un prim pas, vectorul mediu din toate realizările particulare ale vectorului aleator  $x$ . Astfel obținem:

$$a_0^1 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \quad a_0^2 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad a_0^3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad a_0^4 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- a. Prima metodă de estimare a matricei de covarianță se bazează pe utilizarea relației (5.234).

$$\begin{aligned} \hat{C}_x &= E\{(x - m_x)(x - m_x)^{*T}\} \cong \frac{1}{4} \sum_{i=1}^4 (a^i - \hat{m}_x)(a^i - \hat{m}_x)^{*T} = \\ &= \frac{1}{4} \left\{ \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -1 & 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} \begin{bmatrix} 0 & -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \end{bmatrix} \right\} = \begin{bmatrix} 1/2 & 1/4 \\ 1/4 & 1/2 \end{bmatrix} \\ \hat{C}_x &= \begin{bmatrix} 1/2 & 1/4 \\ 1/4 & 1/2 \end{bmatrix} \end{aligned}$$

- b. În cea de a doua metodă de estimare, bazată pe matricea setului de date, construim, mai întâi, matricea  $X_0$  dată de elementele eșantionului din care am extras, anterior, vectorul mediu. Astfel:

$$X_0 = \begin{bmatrix} - & (a_0^1)^{*T} & - \\ - & (a_0^2)^{*T} & - \\ & \vdots & \\ - & (a_0^4)^{*T} & - \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}$$

Matricea de covarianță estimată va fi egală cu:

$$\hat{C}_x = \frac{1}{K} X_0^{*T} X_0 = \frac{1}{4} \begin{bmatrix} -1 & 0 & 1 & 0 \\ 0 & -1 & 1 & 0 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{2} \end{bmatrix}$$

#### 5.5.4. Coeficientul de corelație

**Coeficientul de corelație** este o valoare ce **caracterizează gradul de dependență liniară** dintre două variabile aleatoare și el ia întotdeauna valori în intervalul  $[-1, 1]$ . Cu cât coeficientul de corelație este mai apropiat de capetele intervalului de variație (respectiv, de valorile -1 sau 1), cu atât cele două variabile aleatoare sunt mai puternic dependente liniar una față de cealaltă.

Spunem că avem o relație de *dependență liniară totală*<sup>53</sup> între, să zicem, variabilele aleatoare  $x_k$  și  $x_l$  (reprezentând componentele  $k$  și  $l$  ale unui vector aleator de trăsături), atunci când coeficientul de corelație dintre acestea va avea valoarea 1 sau -1. În această situație particulară, reprezentarea grafică, în spațiul trăsăturilor  $(x_k, x_l)$ , a vectorilor de trăsături, va conduce la un set de puncte poziționate pe o dreaptă cu pantă pozitivă (vezi **Figura 5.23.(f)**) pentru o valoare egală cu +1 a coeficientului de corelație. În cazul în care coeficientul de corelație ia valoarea -1 avem de-a face, din nou, cu o relație de liniaritate perfectă între cele două variabile, iar vectorii de trăsături se vor poziționa pe o dreaptă care va avea de această dată, însă, o pantă negativă.

**Corelația** dintre două variabile aleatoare va avea o **valoare pozitivă** dacă, în medie, creșterea valorii unei variabile aleatoare față de media ei va determina creșterea valorii celeilalte variabile față de propria medie, iar scăderea valorii primei variabile aleatoare față de medie va determina, de cele

53 *Dependența liniară totală* poate fi tradusă, de exemplu, prin aceea că **întotdeauna** creșterea/scăderea valorii unei trăsături față de media ei va determina creșterea/scăderea, cu aceeași cantitate echivalentă, și a valorii celei de a doua trăsături față de media corespunzătoare ei.

mai multe ori, scăderea valorii celei de a doua variabile aleatoare față de propria ei medie. Spunem despre aceste variabile aleatoare că sunt **corelate pozitiv**. Exemple de variabile corelate pozitiv: educația și venitul obținut, consumul grăsimilor saturate și nivelul colesterolului, nivelul colesterolului și probabilitatea unui atac de cord etc.

Pentru variabilele **corelate negativ** mecanismul de variabilitate presupune următoarele: de cele mai multe ori (în medie) creșterea valorii unei variabile aleatoare față de media ei va determina scăderea valorii celeilalte variabile față de propria medie, iar scăderea valorii primei variabile aleatoare față de medie va determina, de cele mai multe ori (în medie), creșterea valorii celei de a doua variabile aleatoare față de propria ei medie. Exemple de variabile corelate negativ: numărul de țigări fumate de mamă și greutatea copilului la naștere, masa mașinii și consumul de carburant pe 100 de Km, numărul de țigări fumate și durata de viață etc.

Dacă coeficientul de corelație ia valoarea zero atunci spunem că **nu există nici o relație linară** între cele două variabile aleatoare.

În prezent se cunosc mai multe metode de calcul al coeficientului de corelație. Unele dintre aceste metode sunt mai potrivite comparativ cu altele, funcție de tipul variabilelor aleatoare analizate. Astfel, spre exemplu, există *coeficientul de corelație Pearson*, *Spearman*, *coeficientul de corelație intraclasă (intraclass correlation coefficient, ICC)* etc.

În cadrul acestei cărți vom prezenta doar **coeficientul de corelație Pearson** care este definit de relația:

$$\rho = \stackrel{def}{\frac{c_{kl}}{\sqrt{c_{kk} \cdot c_{ll}}}} = \frac{Cov(x_k, x_l)}{\sqrt{Cov(x_k)Cov(x_l)}} \quad (5.245)$$

sau, echivalent,

$$\rho = \frac{Cov(x_k, x_l)}{\sqrt{\sigma_{x_k}^2 \cdot \sigma_{x_l}^2}} = \frac{Cov(x_k, x_l)}{\sigma_{x_k} \cdot \sigma_{x_l}} = \frac{E\{(x_k - m_k)(x_l - m_l)^*\}}{\sqrt{E\{(x_k - m_k)^2\} \cdot E\{(x_l - m_l)^2\}}} \quad (5.246)$$

**Problema 5.33:** Pentru următoarele două matrici de covarianță:

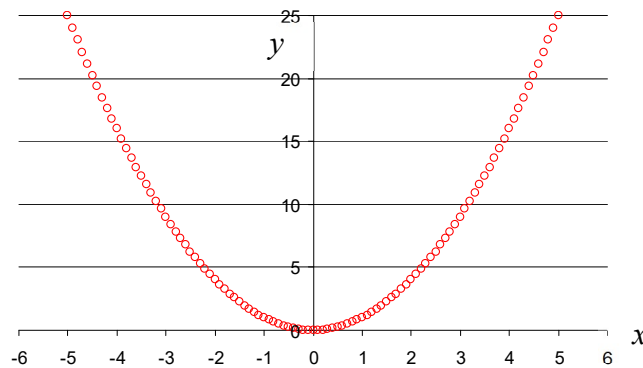
$$C_x = \begin{bmatrix} 0.5 & 0.25 \\ 0.25 & 0.5 \end{bmatrix} \text{ și } C_y = \begin{bmatrix} 4 & -2 \\ -2 & 9 \end{bmatrix}$$

ce caracterizează doi vectori aleatori bidimensionali,  $x$  și  $y$ , determinați coeficientul de corelație între trăsăturile acestor vectori.



**Aplicația 5.4:** Cu ajutorul programului ce se află în directorul **CoefCorr-Pearson**, asociat acestui capitol, și care determină coeficientul de corelație Pearson între două semnale reprezentate vectorial prin doi vectori aleatori  $x$  și  $y$  (indiferent de dimensiunile acestora), se cere să:

1. Determinați relațiile matematice ce ar fi necesare în implementarea programului.
2. Generați cu ajutorul programului ce se află în directorul **Generare distribuție bivariată Gauss-iana** (vezi **Capitolul 6** al acestei cărți) diferite distribuții *Gauss*-iene bidimensionale, al căror coeficient de corelație între variabilele aleatoare ce formează vectorul aleator îl cunoașteți. Converteți seriile de timp astfel încât să poată fi încărcate în programul ce determină coeficientul de corelație – folosiți-vă, de exemplu, de programul de calcul tabelar Microsoft Office Excel. Verificați corectitudinea coeficientului de corelație determinat de program față de cel introdus de dumneavoastră.
3. În directorul seturilor de date asociat acestui capitol se găsesc două serii de timp „*x - distrib parabolica.txt*” și „*y - distrib parabolica.txt*”. În primul fișier avem o variabilă aleatoare uniform distribuită în intervalul  $[-5, 5]$ . În cel de al doilea fișier este un alt șir de valori obținut din primul cu ajutorul relației  $y = x^2$ ; deci, variabila aleatoare  $y$  este complet determinată de variabila aleatoare  $x$ . Dacă le reprezentăm grafic, prima versus cea de a doua, obținem reprezentarea grafică din **Figura 5.28**.



**Figura 5.28.** Reprezentarea grafică a seriilor de timp „*x - distrib parabolica.txt*” și „*y - distrib parabolica.txt*”

Calculând cu ajutorul programului de mai sus coeficientul de corelație vom obține pentru aceste două serii de timp o valoare egală

cu zero, valoare identică cu cea obținută pentru o distribuție de tipul celei prezentate în **Figura 5.23(d)**, în care variabilele aleatoare sunt necorelate. Prin urmare, variabilele aleatoare  $x$  și  $y$  sunt dependente (cunoaștem chiar relația funcțională ce le leagă) dar corelația lor este zero. Din ce cauză am obținut această valoare?

**Rezolvare:**

1. Relația (5.246) se generalizează foarte ușor pentru cazul unor variabile aleatoare  $x$  și  $y$  reale sub forma:

$$\rho = \frac{\text{Cov}(x, y)}{\sqrt{\sigma_x^2 \cdot \sigma_y^2}} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{E\{(x - m_x)(y - m_y)\}}{\sigma_x \cdot \sigma_y} \quad (5.247)$$

Ținând cont, în plus, că:

$$m_x = E\{x\} \quad (5.248)$$

$$\sigma_x^2 = E\{x^2\} - E^2\{x\} \quad (5.249)$$

putem foarte ușor să deducem relația:

$$\rho = \frac{E\{x \cdot y\} - E\{x\}E\{y\}}{\sqrt{E\{x^2\} - E^2\{x\}} \sqrt{E\{y^2\} - E^2\{y\}}} \quad (5.250)$$

sau

$$\rho = \frac{E\{x \cdot y\} - m_x m_y}{\sqrt{E\{x^2\} - m_x^2} \sqrt{E\{y^2\} - m_y^2}} \quad (5.251)$$

Estimând termenii din relația (5.251) conform relației (5.231) și știind că dispunem de un set de date de lungime  $N$  elemente, rezultă imediat că:

$$\rho = \frac{\frac{1}{N} \cdot \sum_{i=1}^N (x^i \cdot y^i) - \frac{1}{N} \sum_{i=1}^N x^i \cdot \frac{1}{N} \sum_{j=1}^N y^j}{\sqrt{\frac{1}{N} \cdot \sum_{i=1}^N (x^i)^2 - \left(\frac{1}{N} \sum_{i=1}^N x^i\right)^2} \sqrt{\frac{1}{N} \cdot \sum_{i=1}^N (y^i)^2 - \left(\frac{1}{N} \sum_{i=1}^N y^i\right)^2}} \quad (5.252)$$

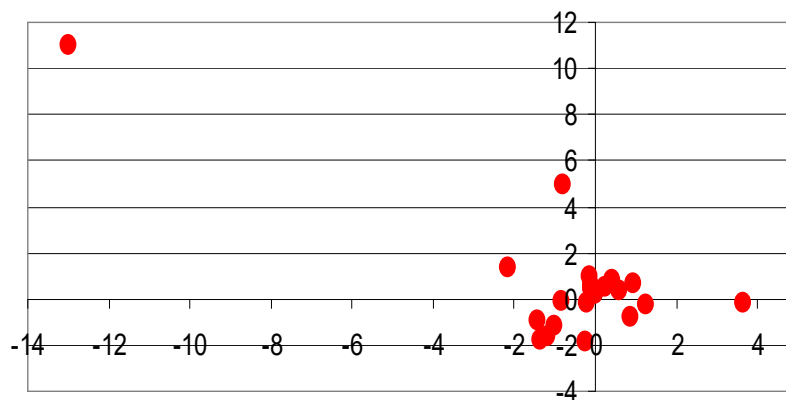
sau

$$\rho = \frac{N \cdot \sum_{i=1}^N (x^i \cdot y^i) - \sum_{i=1}^N x^i \cdot \sum_{j=1}^N y^j}{\sqrt{N \cdot \sum_{i=1}^N (x^i)^2 - \left(\sum_{i=1}^N x^i\right)^2} \sqrt{N \cdot \sum_{i=1}^N (y^i)^2 - \left(\sum_{i=1}^N y^i\right)^2}} \quad (5.253)$$

2. În mod intuitiv pentru distribuția din **Figura 5.28** obținem un coeficient de corelație egal cu zero deoarece valorile mari ale variabilei aleatoare  $y$  sunt asociate atât cu valori mari (pozitive) dar și cu cele mici (negative) ale variabilei aleatoare  $x$ .

*Existența unei relații neliniare* între variabilele aleatoare pentru care calculăm coeficientul de corelație *este o sursă de erori* a coeficientului de corelație Pearson. Acest fapt apare deoarece **coeficientul de corelație Pearson măsoară doar existența sau nu a unei relații liniare** între variabilele aleatoare implicate în calcul.

Anterior am făcut afirmația că dacă *doi vectori aleatori sau variabile aleatoare,  $x$  și  $y$ , sunt independente atunci ele sunt și necorelate*; în schimb, s-a afirmat că *reciproca nu este întotdeauna adevărată*. Inversa afirmației precedente nu este adevărată întotdeauna deoarece corelația pune în evidență doar relațiile liniare ce există între variabile. Dacă parcurgem punctul trei al aplicației precedente observăm existența a două variabile aleatoare necorelate (coeficientul de corelație este zero) dar dependente între ele – variabila aleatoare  $y$  este complet determinată de variabila aleatoare  $x$ , prin intermediul relației  $y = x^2$ .



**Figura 5.29.** Influența elementelor aberante asupra valorii coeficientului de corelație

Atunci când în setul de date apar înregistrări aberante (“*outlier*” în limba engleză), acestea influențează în mod negativ valoarea coeficientului de corelație. Astfel, spre exemplu, dacă analizăm numai setul de date poziționat în centrul sistemului de coordonate, vezi **Figura 5.29**, obținem o valoare a coeficientului de corelație de 0.002853; ținând cont că este un set de date mic, de doar 20 de vectori de trăsături, putem trage concluzia că aceste trăsături, care formează vectorul de trăsături, sunt necorelate. Dacă în schimb, în analiză luăm în calcul și elementul “*aberant*” (cel poziționat în partea stânga sus), atunci obținem un coeficient de corelație al setului de date de -0.7847, valoare pe baza căreia suntem tentați să tragem concluzia că variabilele aleatoare ce caracterizează această distribuție sunt corelate negativ, ceea ce, din nou, nu este adevărat. Raportarea acestui coeficient de corelație, de -0.7847, este falsă deoarece coeficientul astfel calculat nu este reprezentativ pentru cea mai mare parte a setului de date.

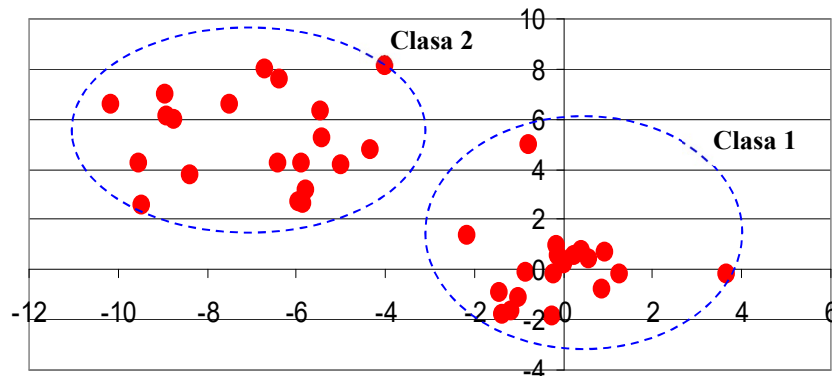
**Aplicația 5.5:** Cu ajutorul programului ce se află în directorul **CoefCorr-Pearson** încărcați seriile de timp „*outliers - x.txt*” și „*outliers - y.txt*”. Vizualizați rezultatele în două situații distincte, respectiv cu și fără primul element din ambele serii. Determinați coeficientul de corelație și reprezentați grafic rezultatele pentru ambele cazuri prezentate precedent, utilizând, de această dată, programul de calcul tabelar Microsoft Office Excel.

Elementele aberante pot să apară, de exemplu, în practică atunci când punctul zecimal nu este poziționat corect (setul de date lucrează cu caracterul „.” drept punct zecimal iar programul ce realizează analiza utilizează drept punct zecimal caracterul „,”) sau atunci când unul din electrozii cu ajutorul căruia s-a făcut achiziția unui semnal biomedical s-a desprins ș.a.m.d.

Datorită erorilor pe care le pot genera aceste elemente aberante, anterior oricărei analize, fiecare variabilă aleatoare este preprocesată și sunt menținuți spre analiză și procesare, de exemplu, numai acele valori ale vectorului aleator de trăsături pentru care valorile elementelor componente (realizări particulare ale variabilelor aleatoare componente) se găsesc, fiecare în parte, în intervalul  $[m - 3\sigma, m + 3\sigma]$  (prin  $m$  am notat valoarea medie a variabilei aleatoare iar  $\sigma$  este deviația standard a aceleiași variabile aleatoare); pentru o distribuție normală în intervalul astfel specificat ar trebui să se regăsească, teoretic cel puțin, cca. 99.72% dintre valori.

O altă problemă ce poate să apară și să influențeze în mod negativ calcularea coeficientului de corelație este dată de existența, în setul de date, a unui număr de clase distincte care, însă, sunt considerate în cadrul analizei,

drept o singură clasă. Astfel, dacă analizăm distribuția prezentată în **Figura 5.30**, ambele clase sunt caracterizate de vectori de trăsături necorelate. Coeficientul de corelație pentru prima clasă este de 0.0028 iar pentru cea de a doua clasă este de -0.0099. Dacă, în schimb, nu am fi dispus de nici o informație *a priori* privind structura setului de date și, în consecință, l-am fi privit ca pe un întreg, ca pe un set de date ce caracterizează o singură clasă, atunci am fi obținut un coeficient de corelație global de -0.773, ceea ce nu ar fi reflectat situația reală existentă.



**Figura 5.30.** Influența existenței a două clase distincte în setul de date asupra coeficientului de corelație

**Aplicația 5.6:** Cu ajutorul programului ce se află în directorul *CoefCorr-Pearson* încărcați seriile de timp „2 clase - x.txt” și „2 clase - y.txt”. Vizualizați și analizați rezultatele în următoarele două situații distincte: **(a)** ambele clase considerate drept un întreg (drept o singură clasă) și **(b)** cazul în care clasele sunt considerate și analizate separat. Determinați coeficientul de corelație și reprezentați grafic rezultatele pentru toate situațiile prezentate anterior, utilizând pentru aceasta programul de calcul tabelar, Microsoft Office Excel.

Din ultimele trei exemple s-a putut observa că există un număr de situații (dependență neliniară a datelor, valori aberante în setul de date, mai multe clase distincte în setul de date etc.) când valorile coeficientului de corelație nu reflectă situația reală și mecanismele de cuplare ce au generat datele respective. Din acest motiv, și nu numai, apare ca fiind **imperios necesar să generăm, într-o primă etapă, o reprezentare vizuală a seturilor de date** și abia ulterior să determinăm valoarea coeficientului de corelație, în mod specific, pentru fiecare set de date.

O aplicație practică, directă, a coeficientului de corelație ține de *reducerea dimensionalității spațiului de trăsături*. Dacă, de exemplu, două trăsături ale unui vector de trăsături,  $d$ -dimensional, sunt caracterizate de un coeficient de corelație apropiat de valorile 1 sau -1, aceste trăsături pot fi combinate într-o singură trăsătură sau, eventual, una dintre ele poate fi eliminată complet. În acest mod, dimensionalitatea vectorului de trăsături se reduce cu o unitate și, simultan, informația redundantă existentă în setul de date este eliminată.

### 5.5.5. Funcția de coerență spectrală

Dependența dintre două semnale staționare (vezi **Anexa: Noțiunea de semnal. Semnal staționar**) și mărginite poate fi caracterizată de anumiți parametri care evaluează și evidențiază corelația dintre acestea.

Pentru două procese  $x(t)$  și  $y(t)$  **funcția de coerență spectrală** este dată de relația:

$$\left| R_{xy}(\omega) \right|^2 = \lim_{T \rightarrow \infty} \left| \text{corr} \left\{ F_x^T(\omega) \cdot F_y^T(\omega) \right\} \right| \quad (5.254)$$

Această funcție furnizează o măsură a gradului de corelație între procesele  $x(t)$  și  $y(t)$  pentru fiecare componentă spectrală în parte. Relația (5.254) poate fi văzută ca pătratul amplitudinii funcției de corelație între transformatele Fourier ale secvențelor  $x(t)$  și  $y(t)$  [Amjad, 1997].

În relația (5.254),  $F_x^T(\omega)$  reprezintă transformata Fourier a unui segment de semnal de lungime  $T$ , aparținând procesului  $x(t)$ .

Deoarece rezultatele obținute în mod practic reprezintă doar niște estimări ale funcției de coerență spectrală, se creează și, în acest caz, premisele calcului estimatorilor iar, în final, prezentăm inclusiv gradul corespunzător de încredere al estimării.

Pentru a calcula un estimat al funcției de coerență, împărțim, mai întâi, seriile de timp în secvențe de o anumită durată/lungime (fie  $T$  lungimea unor astfel de secvențe). În consecință, pentru fiecare serie de timp înregistrată se obțin, să zicem,  $L$  astfel de segmente, fiecare de lungime  $T$ . Transformata Fourier a unui astfel de segment,  $F_x^T(\omega, l)$ ,  $l = 1 \dots L$ , pentru pulsația  $\omega$  este dată de relația:

$$F_x^T(\omega, l) = \int_{(l-1)T \Delta t_s}^{lT \Delta t_s} x(t) \cdot e^{-i\omega t} dt \approx \sum_{t=(l-1)T}^{lT} e^{-i\omega t} x(t) \quad (5.255)$$

Definiția funcției de corelație între transformatele Fourier a celor două procese în funcție de varianța și covarianța acestora pe un anumit segment este dată de:

$$\text{corr} \left\{ F_x^T(\omega) F_y^T(\omega) \right\} = \frac{\text{cov} \left\{ F_x^T(\omega) F_y^T(\omega) \right\}}{\sqrt{\text{var} \left\{ F_x^T(\omega) \right\} \cdot \text{var} \left\{ F_y^T(\omega) \right\}}} \quad (5.256)$$

Folosindu-ne de relația (5.256) putem obține o definiție alternativă a funcției de coerență spectrală între două procesele  $x(t)$  și  $y(t)$ :

$$|R_{xy}(\omega)|^2 = \frac{|f_{xy}(\omega)|^2}{f_{xx}(\omega)f_{yy}(\omega)} \quad (5.257)$$

unde  $f_{xy}(\omega)$  este **densitatea cross-spectrală** a celor două procese  $x(t)$  și  $y(t)$  iar  $f_{xx}(\omega)$  și  $f_{yy}(\omega)$  sunt **densitățile auto-spectrale**, pentru pulsația  $\omega$ , ale fiecărui proces în parte.

Relația (5.257) indică gradul de corelație liniară în domeniul spectral între două semnale, pe o scală de la 0 (semnificând independența liniară) la unu (completa dependență liniară a unui proces de celălalt), acest parametru putând fi interpretat și ca o măsură a predictibilității liniare.

Un avantaj al acestei metode este dat de faptul că, indiferent de unitățile de măsură ale seriilor  $x(t)$  și  $y(t)$ , valoarea coerenței spectrale pentru o anumită frecvență nu are nici o unitate de măsură, furnizând doar o interpretare a relației în care se găsesc cele două semnale, la frecvența respectivă, conform cu specificațiile anterioare.

Auto-spectrul și cross-spectrul necesare în relația (5.257) pot fi estimate prin medierea segmentelor de semnal disjuncte, anterior obținute. Pentru calcularea (estimarea) cross-spectrului relația de calcul este dată de:

$$\hat{f}_{xy}(\omega) = \frac{1}{2\pi LT} \sum_{l=1}^L F_x^T(\omega, l) \cdot \overline{F_y^T(\omega, l)} \quad (5.258)$$

În calcularea auto-spectrului secvenței  $x$ ,  $f_{xx}(\omega)$ , se folosește relația (5.258) în care  $y$  se înlocuiește cu  $x$ . În relația (5.258) bara orizontală de deasupra termenului  $F_y^T(\omega, l)$  simbolizează complex conjugatul acestuia.

În final, pentru estimarea funcției de coerență spectrală din relația (5.257) am utilizat estimatorii calculați cu relația (5.258). Astfel:

$$|\hat{R}_{xy}(\omega)|^2 = \frac{|\hat{f}_{xy}(\omega)|^2}{\hat{f}_{xx}(\omega)\hat{f}_{yy}(\omega)} \quad (5.259)$$