

## 7. Clasificatori statistici

### 7.1. Clasificatorul bazat pe metrica Mahalanobis

Parte din limitările specifice clasificatorului de tip distanță minimă, ce utilizează metrica euclidiană, pot fi înlăturate utilizându-se pentru aceasta o altă metrică, respectiv, metrica Mahalanobis. În acest mod noul clasificator, bazat pe metrica Mahalanobis, va depăși problemele generate de scalarea setului de date și/sau cele generate de existența trăsăturilor corelate.

#### 7.1.1. Distanța standardizată

##### 1. Caracterizarea statistică a unei clase

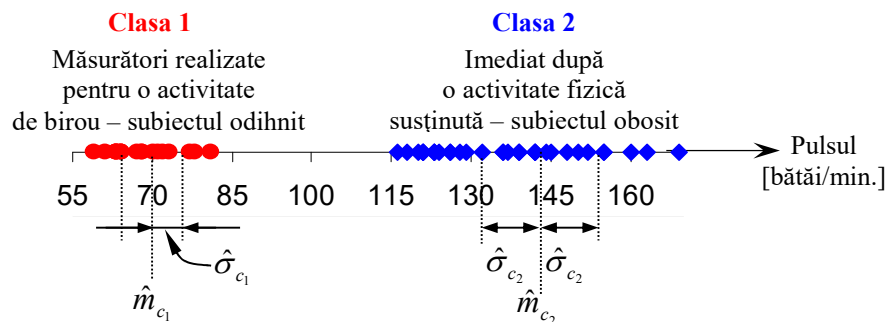
Să considerăm procesul aleator  $x$  caracterizat de un set de  $N$  realizări particulare ale acestuia:  $a^1, a^2, \dots, a^N$ . În continuare vom considera că toate aceste realizări particulare aparțin aceleiași clase. După cum s-a arătat anterior variabila aleatoare  $x$  poate fi caracterizată de doi parametri statistici fundamentali: media și varianța. Media realizărilor particulare ale variabilei aleatoare  $x$  se estimează cu ajutorul relației (5.229) sau a relației (5.230) și poate fi scrisă:

$$\hat{m}_x = \frac{1}{N} (a^1 + a^2 + \dots + a^N) \quad (7.1)$$

Dacă setul de date,  $a^1, a^2, \dots, a^N$ , aparține aceleiași clase (condiție satisfăcută în cazul nostru) valoarea medie este aproximativ centrul clasei respective, vezi **Figura 7.1**. În această situație spunem că media clasei este o valoare tipică pentru clasa respectivă. În cazul clasificatorului de tip minimă distanță vectorul mediu era considerat vectorul prototip al clasei (pentru  $x$ , variabilă aleatoare, media reprezenta valoarea prototip).

Varianța (momentul centrat de ordin doi al variabilei aleatoare  $x$ ) este și ea o măsură ce caracterizează clasa respectivă însă din punct de vedere al mărimii clasei respective, al extinderii ei spațiale; cu alte cuvinte, varianța

reflectă dispersia clasei sau cât de departe poate fi o valoare particulară a lui  $x$  față de cea tipică, dată de media clasei.



**Figura 7.1.** Media și deviația standard pentru cazul unidimensional a două clase

Valoarea estimată a varianței este, așa după cum știm, dată de media aritmetică, calculată la nivel de eșantion, a pătratelor deviațiilor individuale de la medie:

$$\hat{\sigma}_x^2 = \frac{1}{N} \left[ (a^1 - m_x)^2 + (a^2 - m_x)^2 + \dots + (a^N - m_x)^2 \right] \quad (7.2)$$

Rădăcina pătratică a varianței este deviația standard. În **Figura 5.31** se prezintă corelația ce există între repartiția numărului de elemente și deviația standard a distribuției acestor elemente, dacă variabila aleatoare  $x$  este caracterizată de o funcție densitate de probabilitate de tip *gauss*-iană. În acest caz 68% dintre exemplare vor fi la o distanță de maximum o deviație standard față de medie în timp ce 95% dintre exemplare vor fi la o distanță maximă de două deviații standard.

În reprezentarea grafică din **Figura 7.1** se prezintă, în mod intuitiv, două clase și, respectiv, parametrii statistici de ordin unu și doi ce le caracterizează. Din această figură se observă că media (sau generalizând vectorul mediu) ne furnizează informații ce privesc poziționarea clasei (a centrului ei) în spațiul trăsăturilor, în timp ce deviația standard ne furnizează informații despre “raza” sau dimensiunea clasei.

**Observație 7.1:** Același tip de analiză, folosind statistica descriptivă, poate fi făcută și la nivelul fiecărei variabile aleatoare – componentă a vectorului aleator multidimensional  $x$ . Astfel, dacă  $x_i$  reprezintă variabila aleatoare aferentă trăsăturii  $i$  din vectorul aleator de trăsături  $x$ , atunci putem vorbi de  $m_{ji}$  ca fiind media trăsăturii  $i$  pentru clasa  $j$ , și,

respectiv, de  $\sigma_{ji}$  ca fiind deviația standard pentru aceeași trăsătură  $i$  și aceeași clasă  $j$ .

## 2. Definiția distanței standardizate

În general, orice valoare numerică pentru o anumită trăsătură, componentă a unui vector aleator de trăsături,  $x$ , poate avea o unitate de măsură specifică. În acest caz vorbim, implicit, și de o scală de lucru corespunzătoare. Dacă o astfel de trăsătură este multiplicată cu o cantitate  $q$  (factor de scalare) atunci, atât media cât și deviația standard sunt, la rândul lor, multiplicată cu aceeași cantitate  $q$  – în cazul mediei – și, respectiv, cu  $q^2$  – în cazul varianței; aceste rezultate reprezintă particularizări ale relațiilor (6.12) și (6.15) pentru cazul unidimensional.

În concluzie, atât poziționarea dar și dispersia clasei se modifică ca efect al modificării scalei oricăreia dintre componentele vectorului de trăsături – aspect care poate fi în unele aplicații foarte supărător.

Pentru ușurarea calculelor în unele situații este de dorit o scalare a setului de date astfel încât deviația standard a variabilei aleatoare studiate să fie egală cu unitatea. Această operație este foarte ușor de realizat prin împărțirea variabilei aleatoare  $x$  la deviația ei standard,  $\sigma_x$ . În acest mod, pe lângă obținerea efectului dorit, noua variabilă aleatoare rezultantă este o variabilă aleatoare fără unitate de măsură. Această caracteristică este foarte importantă întrucât, pe de o parte, ea rezolvă neajunsul prezentat mai sus, iar pe de altă parte, ne permite definirea unei distanțe care este, de asemenea, una independentă de unitățile de măsură folosite la observarea diverselor trăsături analizate.

Pentru măsurarea distanței de la un element  $a$  la centrul clasei  $m_x$  (ca în cazul clasificatorului de tip minimă distanță) este, deci, util să măsurăm această distanță în mod relativ, prin împărțirea la deviația standard. Această distanță este numită în literatura de specialitate **distanța standardizată** și ea este dată de:

$$r = \left| \frac{a - m_x}{\sigma_x} \right| \quad (7.3)$$

Utilizând relația (7.3) se poate demonstra ușor că distanța  $r$  este invariantă la operațiile de translare și scalare. Aceste observații sugerează o importantă generalizare a clasificatorului de tip minimă distanță bazat pe metrica Euclidiană.

**Problemă 7.1:** Demonstrați faptul că distanța standardizată este invariantă la operațiile de translare și scalare.

În măsurarea distanței dintre un vector oarecare de trăsături  $a$  și vectorul mediu  $m_j$ , reprezentând prototipul clasei  $j$ , putem utiliza, de asemeni, și următoarea distanță standardizată dată de relația:

$$r^2(a, m_j) = \left( \frac{a_1 - m_{j1}}{\sigma_{j1}} \right)^2 + \left( \frac{a_2 - m_{j2}}{\sigma_{j2}} \right)^2 + \dots + \left( \frac{a_d - m_{jd}}{\sigma_{jd}} \right)^2 \quad (7.4)$$

Această distanță are și ea, la rândul ei, aceeași importantă proprietate și anume aceea de a fi invariantă la operațiile de scalare și translare. Astfel, dacă utilizăm această distanță în cadrul unui clasificator, unitățile de măsură pe care noi le utilizăm pentru exprimarea valorilor diferitelor trăsături nu vor mai afecta în nici un mod distanța măsurată între vectori și, în consecință, nu vor mai influența nici rezultatul final al clasificării.

O generalizare directă a acestei distanțe standard pentru cazul multidimensional este și distanța bazată pe metrica *Mahalanobis*, metrică pe care o prezentăm în cele ce urmează.

### 7.1.2. Metrica Mahalanobis

Să presupunem că avem un vector aleator  $d$ -dimensional  $x$  care este caracterizat de media  $m_x$  și de matricea de covarianță  $C_x$ . În continuare utilizăm o matrice  $A$ ,  $d \times d$  dimensională, pentru transformarea vectorului aleator  $x$  într-un vector aleator  $y$  prin intermediul relației:

$$y = A \cdot x \quad (7.5)$$

În cadrul acestui subcapitol dorim să găsim o metrică capabilă de a generaliza conceptul de distanță, astfel încât distanța de la o realizare particulară  $a$  a vectorului aleator  $x$  la  $m_x$  (media clasei) cât și distanța de la  $b$  la  $m_y$  (atât  $b$  cât și  $m_y$  au fost obținuți din  $a$  și, respectiv,  $m_x$  prin intermediul transformatei  $A$ ) să fie egale în ambele spații în care realizăm măsurătorile.

Cu alte cuvinte, având o astfel de metrică nu vom mai fi obligați, în cazul trăsăturilor corelate, să găsim o transformare care să ne ducă din spațiul inițial existent într-un alt spațiu în care matricea de covarianță să fie una diagonală. Cu ajutorul acestei metrici vom putea realiza clasificarea în chiar spațiul inițial, deoarece performanțele metricii vor fi aceleași indiferent de spațiul în care lucrăm (astfel spus, ne este indiferent spațiul unde facem clasificarea). În această nouă abordare putem economisi putere de calcul

întrucât nu mai suntem obligați să aflăm transformata care să determine decorrelarea trăsăturilor și pe care, ulterior, să o mai și aplicăm întregului set de date.

**Observația 7.2:** În cazul utilizării metricii euclidiene distanțe egale pentru o realizare particulară  $a$  în cele două spații (în cel inițial și, respectiv, în cel obținut în urma aplicării transformării) sunt garantate numai în acele cazuri particulare (de exemplu, în situația în care matricea  $A$  doar reflectă întregul spațiu în raport cu un hiperplan de referință) Ceea ce dorim să facem în mod real este să normalizăm distanțele, similar cu situația monodimensională dată de relația (7.3) astfel încât noua distanță obținută să fie invariantă la aplicarea oricărui tip de operator liniar.

Scopul principal al acestui subcapitol este de a generaliza distanța standard, dată de (7.3), la o distanță capabilă să lucreze cu vectori multidimensionali de trăsături.

Dacă rescriem relația (7.3) în forma:

$$r^2 = \left( \frac{a - m_x}{\sigma_x} \right)^2 = (a - m_x) \frac{1}{\sigma_x^2} (a - m_x) \quad (7.6)$$

atunci, o generalizare imediată a distanței unidimensionale standardizate pentru cazul unui vector aleator multidimensional va fi dată de relația:

$$r^2 = (a - m_x)^T \frac{1}{C_x} (a - m_x) \quad (7.7)$$

Distanța anterioară (de la vectorul aleator  $a$  la vectorul  $m_x$ ) este cunoscută sub numele de **distanță Mahalanobis**.

**Problemă 7.2:** Să se demonstreze că distanța Mahalanobis este invariantă la orice transformare liniară aplicată setului de date.

**Rezolvare:**

Să considerăm un vector aleatoriu real  $x$  (caracterizat de media  $m_x$  și matricea de covarianță  $C_x$ ) căruia îi aplicăm transformarea dată de relația (7.5). În final va rezulta un vector aleator  $y$  pentru care avem:

$$m_y = A m_x \quad (7.8)$$

și

$$C_y = A C_x A^T \quad (7.9)$$

Știind că:

$$C_y^{-1} = (A^{-1})^T C_x^{-1} A^{-1} \quad (7.10)$$

putem scrie:

$$\begin{aligned} r_y^2 &= (y - m_y)^T C_y^{-1} (y - m_y) = \\ &= (Ax - Am_x)^T (A^{-1})^T C_x^{-1} A^{-1} (Ax - Am_x) = \\ &= (x - m_x)^T A^T (A^{-1})^T C_x^{-1} A^{-1} A (x - m_x) = \\ &= (x - m_x)^T (A^{-1} A)^T C_x^{-1} (A^{-1} A) (x - m_x) = \\ &= (x - m_x)^T C_x^{-1} (x - m_x) = \\ &= r_x^2 \end{aligned} \quad (7.11)$$

Această demonstrație dovedește invarianța distanței Mahalanobis la orice tip de transformare liniară.

Se poate arăta că suprafețele de decizie date de norma Mahalanobis sunt pătratice iar pentru distanțe  $r$  constante față de  $m_x$  obținem elipsoizi care sunt centrați în acest vector.

### 7.1.3. Clasificatorul Mahalanobis

Un clasificator de minimă distanță ce utilizează distanța Mahalanobis se numește **clasificator Mahalanobis**. În continuare prezentăm modul de utilizare a distanței Mahalanobis în cadrul unui clasificator de tip minimă distanță.

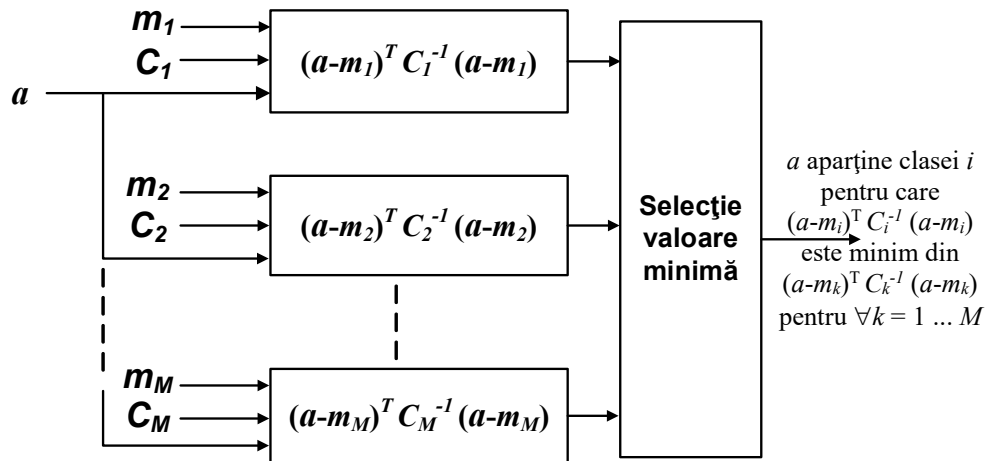


Figura 7.2. Schema bloc a clasificatorului bazat pe norma Mahalanobis

Fie  $m_1, m_2, \dots, m_M$  centrele de masă (șabloanele, *template*-urile) pentru cele  $M$  clase considerate (implicit se presupune că acestea formează o partiție a eșantionului), iar  $C_1, C_2, \dots, C_M$  matricile de covarianță corespunzătoare. Cu ajutorul normei Mahalanobis se măsoară distanțele de la vectorul  $a$  la fiecare element prototip aparținând celor  $M$  clase. Vectorul  $a$  va fi, în final, atribuit acelei clase pentru care distanța Mahalanobis de la el la media clasei este minimă. Schema bloc a acestui clasificator este prezentată în **Figura 7.2**.

**Observația 7.3:** În cazul particular când trăsăturile sunt necorelate (matricea de covarianță,  $C_x$ , este diagonală) iar varianțele în toate direcțiile spațiului de trăsături sunt aceleași (toate elementele de pe diagonala principală a lui  $C_x$  sunt egale între ele), suprafețele de decizie devin hipersfere iar distanța Mahalanobis, în această situație, este identică cu distanța Euclidiană. În concluzie, un clasificator de minimă distanță ce utilizează distanța Mahalanobis este, de fapt, o **generalizare a unui clasificator de tip minimă distanță ce utilizează distanța Euclidiană**.

#### 7.1.4. Avantaje și dezavantaje ale clasificatorului Mahalanobis

Utilizarea metricii Mahalanobis înlătură o serie de limitări prezente în cazul utilizării metricii Euclidiene, asigurând totodată următoarele avantaje:

- (1) independență față de scalarea uneia sau alteia dintre trăsături;
- (2) capacitate de a lucra cu trăsături corelate;
- (3) flexibilitate superioară a suprafețelor de decizie, care de această dată sunt unele pătratice.

Toate aceste avantaje prezentate anterior au și un preț pe care trebuie să-l plătim. Astfel:

- (1) Pentru utilizarea metricii Mahalanobis trebuie să estimăm (din setul de date) matricile de covarianță pentru fiecare clasă în parte. Din păcate însă aceste matrici sunt extrem de dificil de estimat cu acuratețe. Mai mult,  timpul de calcul cât și memoria de stocare crește pătratic cu numărul de trăsături pe care îl utilizăm.

Dacă avem  $N$  realizări particulare pentru vectorul aleator real  $x, a^1, a^2, \dots, a^N$ , toate aparținând aceleiași clase, matricea de covarianță va fi calculată cu relația:

$$C_x = \frac{1}{N} \left[ (a^1 - m_x)(a^1 - m_x)^T + (a^2 - m_x)(a^2 - m_x)^T + \dots + (a^N - m_x)(a^N - m_x)^T \right] \quad (7.12)$$

Dacă ținem cont că vectorul aleator  $x$  este unul  $d$ -dimensional rezultă că matricea de covarianță  $C_x$  este o matrice  $d \times d$  dimensională.

- (2) Dacă  $N$ , numărul total de trăsături cu ajutorul cărora estimăm  $C_x$ , este mai mic decât  $d + 1$  atunci ne aflăm în situația particulară în care matricea  $C_x$  este o matrice singulară<sup>1</sup>; acest fapt generează o problemă majoră, întrucât noi suntem obligați să inversăm matricea de covarianță pentru a putea calcula mai departe distanța Mahalanobis. Pentru ca o matrice de covarianță să nu fie singulară, la limită, va trebui să avem un număr minim de vectori de trăsături care să fie mai mare de  $d + 1$ . Lucrurile se complică, însă, și mai mult atunci când conștientizăm că numărul de vectori aleatori, cu ajutorul cărora se estimează matricea de covarianță pentru fiecare clasă în parte, trebuie să fie mai mare decât acest prag. În consecință, pentru o aplicație cu  $M$  clase și un vector de trăsături  $d$ -dimensional setul minim de date ar trebui să conțină cel puțin  $M \times (d + 1)$  vectori de trăsături, cu minim  $d + 1$  vectori aleatori pentru fiecare clasă.

**Observația 7.4:** Chiar în situația în care numărul de vectori aleatori ai setului de date va fi puțin mai mare de  $d + 1$ , estimarea matricei de covarianță nu va fi una foarte exactă. Prin utilizarea unei matrice de covarianță a cărei estimare nu este foarte exactă vom obține o acuratețe scăzută atât în cazul clasificatorului de tip minimă distanță ce utilizează o normă Mahalanobis cât și în cazul clasificatorului *Bayes*-ian, așa cum de altfel vom prezenta și în subcapitolul următor. În ceea ce privește existența unui număr limitat de vectori de trăsături aceasta reprezintă un impediment major în cazul ambelor clasificatoare (*Bayes*-ian și, respectiv, a celui bazat pe norma *Mahalanobis*).

- (3) Matricea de covarianță (pentru vectori aleatori  $d$ -dimensionali) conține  $d^2$  elemente, din care – datorită simetriei acestei matrici – doar un număr de  $d(d+1)/2$  elemente sunt independente. În acest caz putem spune că o estimare corectă a matricii  $C_x$  se obține efectiv abia în momentul în care numărul de vectori de trăsături cu ajutorul cărora calculăm această matrice ajunge aproape de  $d(d + 1)/2$  sau, ideal, depășește această valoare. Dacă problema de rezolvat este una „mică”

<sup>1</sup> O matrice este singulară dacă determinantul ei este egal cu zero.



din punct de vedere al numărului trăsăturilor implicate, acest fapt nu ar ridica probleme deosebite. Dar, din păcate, nu este ieșit din comun să avem probleme de clasificare caracterizate de vectori de trăsături ce au, de exemplu, în jur de 100 de componente.

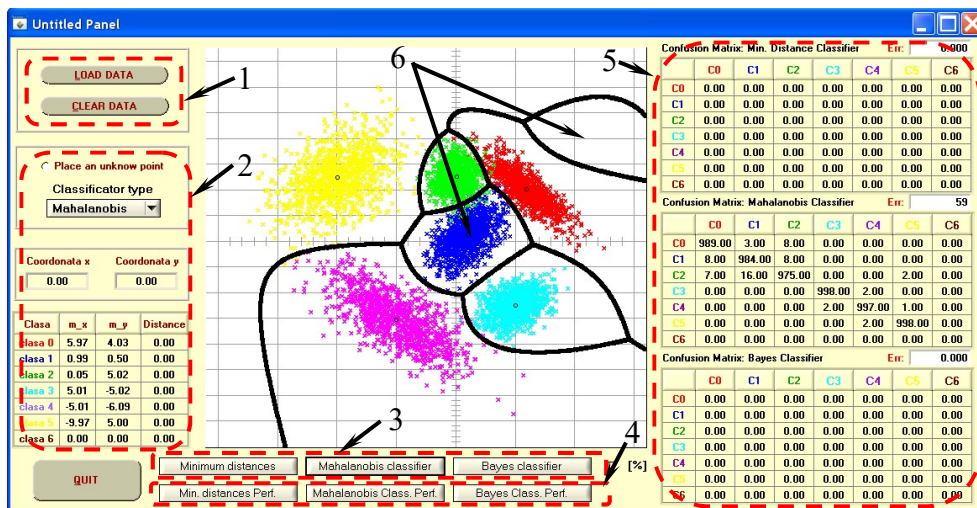
**Exemplu 7.1:** În cazul unei probleme de clasificare a 5 *task*-uri mentale pe baza coeficienților de amplitudine ai modelului ANAPP (*Adaptive Nonlinear Amplitude and Phase Process*) – model prezentat în referința [Dobrea, 2007] – s-au utilizat vectori de trăsături cu un număr de 104 componente; aceste componente reprezentau, concatenat într-un singur vector de trăsături, diverși parametri estimați pe un număr de 6 secvențe de semnal EEG înregistrate simultan de la 6 electrozi plasați la nivelul scalpului unui subiect. În acest caz, o realizare particulară a vectorului aleator de trăsături a fost asociată unor înregistrări (simultane) ale semnalului EEG de la cei 6 electrozi în timp ce setul de date a fost generat de  $N$  astfel de înregistrări simultane. În această aplicație particulară  $d$  a fost egal cu 104 iar un calcul simplu ne relevă un necesar de cel puțin  $d(d+1)/2$ , adică de cel puțin 5460 de realizări particulare pentru o estimare corectă a matricei de covarianță pentru doar o singură clasă. Ținând cont că aplicația de clasificare viza 5 astfel de clase, rezultă că în total ar fi fost nevoie de minim 27300 de vectori de trăsături pentru o estimare corectă a celor cinci matrici de covarianță. În realitate baza de date utilizată a avut doar 1670 de vectori de trăsături [Dobrea, 2007]. Dacă în această situație s-ar fi utilizat un clasificator de tipul Mahalanobis, rezultatele clasificării ar fi fost aproape cu siguranță unele foarte slabe. Motivul unor astfel de performanțe scăzute nu s-ar fi datorat puterii scăzute de discriminare a acestei metode, așa cum am fi fost poate tentați să credem la o primă analiză ci ele ar fi fost generate, în principal, de o estimare nesatisfăcătoare a parametrilor modelului.

În concluzie, se observă că o dată cu creșterea numărului de trăsături apar o serie de factori limitativi ce determină o acuratețe scăzută a clasificatorului bazat pe metrica Mahalanobis.

**Aplicație 7.1:** Pentru înțelegerea clasificatorului Mahalanobis s-a realizat un program al cărui cod se află în directorul „**Comparatie minDist-Mahalanobis-Bayes**” asociat acestui subcapitol. Programul este unul mai general și el implementează trei clasificatori elementari (clasificatorul de minimă distanță bazat pe metrica euclidiană,

clasificatorul *Mahalanobis* și, respectiv, clasificatorul *Bayes-ian*, ce va fi prezentat mai târziu). Programul permite:

- încărcarea unui set de vectori bidimensionali de trăsături, set de date ce poate conține un număr de maxim 7 clase – operația se realizează prin intermediul elementelor 1 de pe interfața grafică (vezi **Figura 7.3**);
- plasarea în spațiul trăsăturilor, cu ajutorul *mouse*-ului, a unui element necunoscut urmând ca, ulterior, să se determine apartenența acestuia la una din clasele posibile utilizând, pentru aceasta, unul din clasificatorii anterior menționați; programul va prezenta simultan și valorile funcțiilor discriminant calculate pentru elementul necunoscut – pentru aceasta se folosesc elementele din registrul 2;
- trasarea, pentru setul de date încărcat, a suprafețelor de decizie obținute folosind cei trei clasificatori – butoanele dedicate sunt cele din registrul 3;
- Afișarea, în valori absolute sau procentuale, a performanțelor clasificării obținute folosind succesiv, pe același set de date, cei trei clasificatori; pentru afișarea performanțelor s-a folosit ca formă de prezentare matricea confuziilor.



**Figura 7.3.** Interfața grafică a programului

Utilizând acest program se cere să se atingă următoarele obiective:

1. Încărcați unul din următoarele seturi de date: „6 Classes.txt”, „7 Classes.txt” (aceste două seturi de date conțin 6, respectiv, 7 clase,

cu număr similar de vectori de trăsături bidimensionali – 1000 pentru fiecare clasă în parte; aceste clase sunt caracterizate, fiecare în parte, de o matrice de covarianță și, respectiv, de o medie proprie), „*6 Classes - Different prob.txt*” (acest set de date este derivat din „*6 Classes.txt*”, de care diferă prin aceea că cea de a treia clasă este compusă din doar 10% din elementele clasei originale), „*Femei-Barb Greut-Inalt Bayes rnd.txt*”, „*Odihnit-Obosit Puls-PresSist rnd.txt*”.

2. Determinați generarea suprafețelor de decizie pentru clasificatorul de minimă distanță și pentru cel *Mahalanobis*.
3. Determinați performanțele de clasificare ale acestor doi clasificatori (analizate prin intermediul matricilor confuziilor) și corelați rezultatele obținute cu poziționarea suprafețelor de decizie anterior obținute.
4. Plasați elemente necunoscute în spațiul de intrare și analizați corecta lor atribuire la una din clasele existente prin intermediul clasificatorului de minimă distanță și a celui de tip *Mahalanobis*.
5. Prin trasarea suprafețelor de decizie pentru clasificatorul *Mahalanobis* și analizând rezultatele obținute la punctul precedent se poate observa că una din clase este caracterizată de existența a două regiuni decizionale disjuncte (respectiv, regiunile 6); dintre aceste regiuni, una nici măcar nu conține elemente. Explicați, în mod intuitiv, ce anume a determinat obținerea acestui rezultat.

## 7.2. Clasificatorul *Bayes-ian*

### 7.2.1 Alegerea optimă a suprafeței de decizie bazată pe modelul statistic al datelor

Clasificatorii prezentați până acum și-au bazat regula de decizie pe rezultatul măsurării unor distanțe (distanțele dintre forma de intrare și un set de vectori de referință sau puncte prototip din spațiul caracteristicilor). O nouă abordare a problemelor de clasificare – bazată de această dată pe modelul statistic al setului de date (vezi definiția modelului statistic dată în *Subcapitolul 5.5.1*) – este prezentată în cele ce urmează.

Ca element de referință, *tehnicele statistice de clasificare* se bazează pe o presupunere fundamentală și anume aceea că pentru fiecare clasă există o funcție densitate de probabilitate ce ne ajută în determinarea probabilității ca o formă de intrare să aparțină unei clase sau alteia dintre clasele posibile.

**Exemplul 7.2:** Determinarea stării de oboseală fizică a unui subiect funcție de activitatea cardiacă a acestuia – cuantizată prin intermediul pulsului persoanei analizate (problemă prezentată anterior, în *Subcapitolul 3.5.1*) –, poate fi modelată, de exemplu, și tratată statistic astfel: dacă considerăm pulsul, măsurat în bătăi/minut, o cantitate guvernată de legile statisticii și generată de două “processe” diferite – oameni odihniți și, respectiv, oameni obosiți<sup>2</sup> fizic – observăm oportunitatea utilizării unui model statistic în cadrul acestei probleme. Cele două clase sunt date, în acest caz, de clasa {odihnit} și, respectiv, clasa {obosit}.

Din măsurări repetate ale pulsului pentru cele două clase de subiecți putem extrage **parametrii statistici** ce determină în mod unic **funcțiile densitate de probabilitate teoretice** propuse pentru a modela distribuțiile celor două clase.

În cazul ipotezei unor distribuții *gauss-iene* monodimensionale – ipoteză folosită și de noi în cazul problemei de față – avem nevoie să estimăm doar doi parametri, respectiv, media și varianța, în timp ce, în cazul distribuțiilor *gauss-iene* caracterizate de vectori *d*-dimensionali de trăsături ar fi trebuit să calculăm vectorul mediu și matricea de covarianță. O dată estimați acești parametri, pentru ambele clase,

<sup>2</sup> În această analiză starea de oboseală fizică a fost indusă de o activitate fizică intensă – urcarea unui deal într-un timp minim posibil fiecărui subiect

putem ulterior **determina complet forma funcțională a densităților de probabilitate** ce le caracterizează. Folosind, în plus, și informațiile *a priori* privind probabilitatea claselor putem deduce mai departe, cu ajutorul regulei lui Bayes, care sunt **probabilitățile posterioare ale claselor** (probabilitățile condiționate revizuite ale claselor, reprezentând probabilitatea ca o realizare particulară să aparțină unei anumite clase). În final, folosind aceste probabilități posterioare ale claselor putem trece la implementarea clasificatorului statistic.

**Observația 7.5:** Problema de mai sus poate fi una generalizată astfel: *avem la dispoziție o realizare particulară,  $a^0$ , a unui vector aleator  $x$ ,  $d$ -dimensional, de trăsături și ne punem problema cărei clase din cele  $M$  posibile să îl atribuim.*

Soluția cea mai logică și directă ar fi să **atribuim această realizare particulară  $a^0$  acelei clase ce are probabilitatea maximă de a-l conține**. În concluzie, observăm că din nou teoria probabilităților și statistica ne poate da o mână de ajutor în rezolvarea acestor clase de probleme.

Teoria probabilităților și statistica ne propun un set de reguli foarte generale și precise (aplicabile la o clasă foarte largă de probleme) pentru construcția unui clasificator. Această teorie ne arată că un clasificator optimal va alege clasa  $c_i$  de apartenență a elementului  $a^0$  (o realizare particulară a unui vector aleator  $x$ ) drept acea clasă ce maximizează probabilitatea condiționată  $P(c_i | a^0)$ . Acest ultim termen,  $P(c_i | a^0)$ , numit și probabilitate posterioară (vezi **Subcapitolul 5.4.5**, „Probabilitatea condiționată. Regula produsului, a sumei și teorema lui Bayes”), reprezintă probabilitatea clasei  $c_i$  de a include elementul  $a^0$ . Din această perspectivă vom avea:

$$a^0 \text{ va fi asignat clasei } c_i \text{ dacă: } P(c_i | a^0) = \max_{j=1 \dots M} P(c_j | a^0) \quad (7.13)$$

În relația (7.13)  $M$  este numărul total de clase căroră elementul  $a^0$  poate să le aparțină virtual; de exemplu pentru problema practică prezentată la începutul acestui subcapitol  $M = 2$  (obosit versus odihnit).

**Observația 7.6:** Din punctul de vedere al teoriei clasificatorilor (vezi **Subcapitolul 3.5.2**) se constată că în relația (7.13) termenii  $P(c_i | a^0)$  sunt chiar funcțiile discriminant caracteristice fiecărei clase.

Deoarece probabilitatea posterioară nu poate fi determinată în mod direct din măsurători asupra setului de date, vom utiliza relația lui *Bayes* pentru a o deduce. Astfel, putem scrie:

$$P(c_i | a^0) = \frac{f_x(a^0 | c_i)P(c_i)}{f_x(a^0)} \quad (7.14)$$

unde:

- $P(c_i | a^0)$  este probabilitatea<sup>3</sup> clasei  $c_i$  de a include elementul  $a^0$  (este probabilitatea ca  $a^0$  chiar să aparțină clasei  $c_i$ );
- $f_x(a^0 | c_i)$  probabilitatea realizării lui  $a^0$  dată de modelul statistic al vectorului aleator  $x$  ce descrie clasa  $c_i$  (valoarea calculată în  $a^0$  a funcției densitate de probabilitate care descrie distribuția vectorului de trăsături pentru clasa  $c_i$ ), mai simplu este probabilitatea ca elementul  $a^0$  să fie generat de clasa  $c_i$ ;
- $P(c_i)$  este probabilitatea apriorică a clasei, calculată în cazul problemei noastre de clasificare ca fiind probabilitatea de realizare a evenimentelor din clasa  $c_i$ ;
- $f_x(a^0)$  este probabilitatea ca evenimentul  $a^0$  să se întâmple indiferent de clasa căreia îi aparține; în general, această probabilitate se calculează cu relația:

$$f_x(a^0) = \sum_{i=1}^M f_x(a^0 | c_i)P(c_i) \quad (7.15)$$

După cum se poate observa din relația (7.14), probabilitatea posterioară poate fi calculată ca produsul dintre probabilitatea apriorică a clasei,  $P(c_i)$ , și probabilitatea ca elementul  $a^0$  să fie generat de un proces caracterizat de clasa  $c_i$ , totul normalizat la  $f_x(a^0)$ .

Determinarea numerică a probabilității posterioare  $P(c_i | a^0)$  și asignarea unui element la o anumită clasă este foarte simplă, ea presupunând parcurgerea următorilor pași:

- (1) ținând cont de setul de date și de presupunerile inițiale făcute asupra formei funcțiilor densitate de probabilitate condiționată,  $f_x(a | c_i)$ , se estimează aceste funcții pentru fiecare clasă în parte;
- (2) se estimează  $P(c_i)$ , probabilitățile apriorice, din datele existente, drept probabilitatea de realizare a evenimentelor din clasa  $c_i$ ;

<sup>3</sup> Ca o observație,  $P(c_i | a^0) = p(c_i | a^0)$ , unde  $p(\cdot)$  este funcția masă de probabilitate condiționată a clasei  $c_i$  dată de  $\{x = a\}$ .

- (3) ulterior, din estimarea funcțiilor densitate de probabilitate ce descriu distribuțiile claselor (realizată în prima etapă), se determină probabilitatea  $f_x(a^0|c_i)$ , reprezentând valoarea funcției  $f_x(a|c_i)$  calculată în punctul  $a^0$ ;
- (4) ultimul factor,  $f_x(a^0)$ , este un factor de normalizare ce este, de obicei, eliminat în cadrul aplicațiilor de clasificare datorită:
- lipsei de informații pe care le aduce în procesul de decizie finală a clasificatorului, și
  - încărcării computaționale inutile.
- (5) În ultimul pas, folosindu-ne de relația (7.14), condiția (7.13) devine:

**$a^0$  va fi asignat clasei  $c_i$  dacă:**

$$P(c_i)f_x(a^0|c_i) = \max_{j=1..M} P(c_j)f_x(a^0|c_j) \quad (7.16)$$

După cum am prezentat anterior, în relația (7.16) termenul  $f_x(a^0)$ , termen independent de clasa  $c_j$ , a fost omis deoarece este același factor normalizator, comun pentru toate clasele.

**Problema 7.3:** Să se determine suprafața optimă de decizie ce separă două clase de subiecți umani (odihniți și oboșiți din punct de vedere fizic), funcție de activitatea cardiacă proprie, utilizând pulsul drept trăsătură discriminatorie între cele două clase. În cadrul acestei probleme se presupune că înregistrările activității cardiace au fost realizate pe un grup de 100 de subiecți (angajați ai aceleiași secții, din aceeași întreprindere) la ora 8.00 AM, imediat după sosirea acestora la serviciu. Așa după cum li s-a explicat și cerut anterior, jumătate din subiecți au parcurs ultima parte a drumului către serviciu pe jos, într-un ritm ceva mai alert iar restul au venit la serviciu cu ajutorul mijloacelor de transport în comun (astfel încât starea de oboseală fizică să nu se instaleze).

Seturile de date pentru un eșantion statistic de 100 respectiv 10000 de înregistrări se găsesc în directorul asociat acestui capitol în fișierele „*Puls Obosit-Odihnit 100.txt*” și, respectiv, „*Puls Obosit-Odihnit 10K.txt*”.

**Rezolvare:** Această problemă va fi rezolvată pentru un eșantion de 100 de înregistrări (o singură zi) și, respectiv, un eșantion de 10000 de înregistrări (efectuate în 100 de zile consecutive), scopul urmărit fiind

acela de a pune în evidență influența mărimii eșantionului statistic<sup>4</sup> asupra rezultatelor obținute.

Dacă particularizăm relațiile anterioare ce descriu clasificatorul *Bayes*-ian, la problema de clasificare a stării de oboseală funcție de activitatea cardiacă a subiecților, în relația (7.16) avem  $i = 1, 2$  și, respectiv  $M = 2$  (două clase de subiecți: odihniți și obosiți).

În ipoteza că funcțiile densitate de probabilitate condiționată,  $f_x(a|c_i)$ , pentru cele două clase sunt densități *gauss*-iene, estimăm parametrii acestora – media și, respectiv, varianța – cu ajutorul relațiilor:

$$\hat{m}_i = \frac{1}{K_i} \sum_{j=1}^{K_i} a_i^j \quad (7.17)$$

$$\hat{\sigma}_i^2 = \frac{1}{K_i} \sum_{j=1}^{K_i} (a_i^j - \hat{m}_i)^2 \quad (7.18)$$

unde  $K_i$  este numărul de realizări particulare ale variabilei aleatoare  $x$  (pulsul subiectului) pentru fiecare clasă  $i$  iar  $a_i^j$  este a  $j$ -a realizare particulară a variabilei aleatoare  $x$  pentru clasa  $i$ .

Parametrii statistici estimați pentru cele două eșantioane sunt cei prezentați în tabelul atașat **Figurii 7.4**, în timp ce o reprezentare grafică a distribuțiilor celor două clase este dată în **Figura 7.4** pentru cazul eșantionului de 10.000 de înregistrări.

O dată cu estimarea parametrilor media și varianța pentru distribuțiile condiționate  $f_x(a|c_i)$  ale celor două clase de indivizi, obținem, practic și o determinare completă a acestor funcții:

$$f_x(a|c_1) = \frac{1}{\sqrt{2\pi\hat{\sigma}_1}} e^{-\frac{1}{2} \left[ \frac{(a-\hat{m}_1)^2}{\hat{\sigma}_1^2} \right]} \quad (7.20)$$

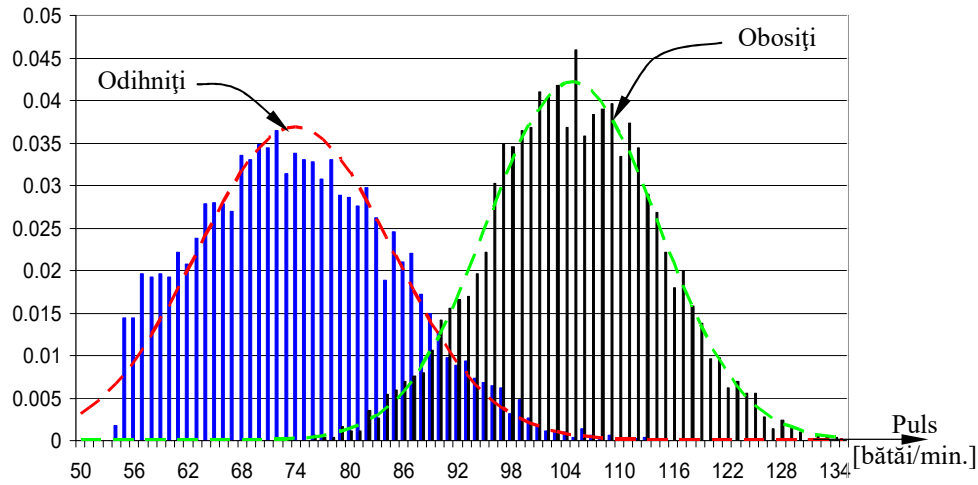
$$f_x(a|c_2) = \frac{1}{\sqrt{2\pi\hat{\sigma}_2}} e^{-\frac{1}{2} \left[ \frac{(a-\hat{m}_2)^2}{\hat{\sigma}_2^2} \right]} \quad (7.21)$$

Probabilitatea  $f_x(a^0|c_i)$ , ca valoarea  $a^0$  să fi fost generată de un subiect aparținând uneia din cele două clase ( $\{\text{odihnit}\}$  sau  $\{\text{obosit}\}$ ),

<sup>4</sup> În cadrul acestei probleme, prin eșantion statistic înțelegem numărul de înregistrări și nu numărul de subiecți care au participat la obținerea setului de date.



se determină prin particularizarea funcțiilor densitate  $f_x(a|c_i)$  în valoarea  $a^0$ .



	10000 de înregistrări	100 de înregistrări
<b>Odihniți</b>	Media = 73.92 Deviația standard = 10.82	Media = 74.32 Deviația standard = 10.88
<b>Obosiți</b>	Media = 104.08 Deviația standard = 9.46	Media = 105.15 Deviația standard = 9.56

**Figura 7.4.** Distribuția seturilor de date: cu linie punctată sunt reprezentate funcțiile densitate de probabilitate ce caracterizează cele 2 clase iar prin segmente verticale sunt reprezentate histogramele celor două clase.

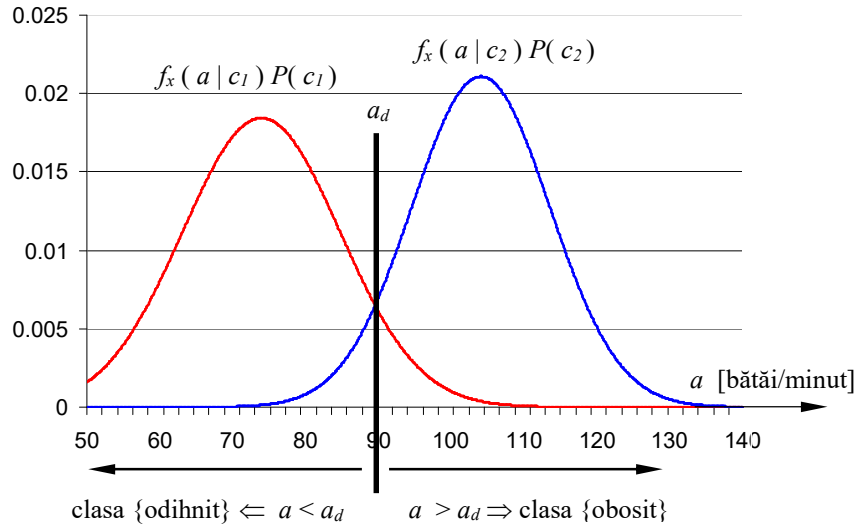
*Termenii  $P(c_i)$* , reprezentând probabilitățile apriorice ale claselor, ar putea fi estimați folosind pentru aceasta, pe lângă informațiile furnizate de setul de date, și informații cu caracter mai subiectiv legate de regimul de viață al subiecților, predispoziția acestora de a efectua, în mod uzual, activități fizice, starea de sănătate<sup>5</sup>, corectitudinea și conștiința angajaților (de exemplu, în a respecta un orar fix), sexul<sup>6</sup>,

<sup>5</sup> Numărul de bătăi pe minut pentru o persoană sănătoasă în stare de repaus este în intervalul [60, 100]. Ritmuri cardiace de peste 100 bătăi/minut (90 după alți autori) se consideră drept patogene – tahicardie. Valori ale activității cardiace sub o valoare de 60 poartă numele de bradicardie. Pentru anumite clase de subiecți (sportivi, adulți tineri și sănătoși) bradicardia este, în anumite condiții, considerată normală.

<sup>6</sup> În medie inima unui bărbat bate în jur de 73 de ori pe minut, iar pentru o femeie ritmul cardiac este în jur de 80 bătăi/minut.

vârsta<sup>7</sup> etc. În general însă, în problemele de clasificare și de învățare, probabilitățile apriorice ale claselor se calculează într-un mod mai obiectiv, ca frecvență relativă a evenimentelor din clasa  $c_i$ :

$$P(c_i) = \frac{K_i}{\sum_{j=1}^M K_j} \quad (7.22)$$



**Figura 7.5.** Determinarea suprafeței de decizie optime,  $a_d$ , utilizând un clasificator *Bayes*-ian

Alegerea poziției suprafeței de decizie  $a = a_d$  (Figura 7.5) este dată de punctul în care probabilitățile posterioare calculate mai sus sunt egale, fiind deci soluția ecuației:

$$P(c_1 | a) = P(c_2 | a) \quad (7.23)$$

de unde rezultă mai departe:

$$P(c_1) \frac{1}{\sqrt{2\pi}\hat{\sigma}_1} e^{-\frac{1}{2}\left(\frac{a-\hat{m}_1}{\hat{\sigma}_1}\right)^2} = P(c_2) \frac{1}{\sqrt{2\pi}\hat{\sigma}_2} e^{-\frac{1}{2}\left(\frac{a-\hat{m}_2}{\hat{\sigma}_2}\right)^2} \quad (7.24)$$

<sup>7</sup> Inima unui nou născut bate de 120 de ori pe minut. Acest ritm scade o dată cu înaintarea în vârstă, la 10 ani ajungând în jur de 90 de bătăi pe minut. Pentru o persoană matură ritmul cardiac mediu este prezentat în observația anterioară.

$$P(c_1)\hat{\sigma}_2 e^{-\frac{1}{2}\left(\frac{a-\hat{m}_1}{\hat{\sigma}_1}\right)^2} = P(c_2)\hat{\sigma}_1 e^{-\frac{1}{2}\left(\frac{a-\hat{m}_2}{\hat{\sigma}_2}\right)^2} \quad (7.25)$$

Logaritmand, obținem:

$$\ln(P(c_1)) + \ln(\hat{\sigma}_2) - \frac{1}{2}\left(\frac{a-\hat{m}_1}{\hat{\sigma}_1}\right)^2 = \ln(P(c_2)) + \ln(\hat{\sigma}_1) - \frac{1}{2}\left(\frac{a-\hat{m}_2}{\hat{\sigma}_2}\right)^2 \quad (7.26)$$

În final rezultă următoarea ecuație de gradul 2:

$$\frac{a^2}{2}\left(\frac{1}{\hat{\sigma}_2^2} - \frac{1}{\hat{\sigma}_1^2}\right) + a\left(\frac{\hat{m}_1}{\hat{\sigma}_1^2} - \frac{\hat{m}_2}{\hat{\sigma}_2^2}\right) + \left[\ln\left(\frac{P(c_1)}{P(c_2)}\right) + \ln\left(\frac{\hat{\sigma}_2}{\hat{\sigma}_1}\right) - \frac{1}{2}\left(\frac{\hat{m}_1^2}{\hat{\sigma}_1^2} - \frac{\hat{m}_2^2}{\hat{\sigma}_2^2}\right)\right] = 0 \quad (7.27)$$

### Discuții:

Pentru varianțe diferite ( $\hat{\sigma}_1 \neq \hat{\sigma}_2$ ) relația (7.27) este o ecuație de gradul doi, cu două soluții. Dacă varianțele  $\hat{\sigma}_1$  și  $\hat{\sigma}_2$  sunt egale atunci relația (7.27) devine o ecuație de gradul întâi, având o singură soluție.

- a) În cazul aplicației noastre, pentru care  $\hat{\sigma}_1 \neq \hat{\sigma}_2$ , ecuația are două rădăcini însă dintre acestea numai una este soluția căutată, respectiv, soluția plauzibilă din punct de vedere biomedical. Soluțiile găsite sunt:
- pentru eșantionul cu 10000 de determinări:  $a_{d1} = 89.55$  și  $a_{d2} = 314.31$  iar
  - pentru eșantionul cu 100 de determinări:  $a_{d1} = 90.29$  și  $a_{d2} = 328.85$ .

Întrucât a doua soluție obținută nu este nici într-un caz, nici în altul, plauzibilă din punct de vedere biomedical, rezultă că suprafața de decizie (care pentru această problemă este un punct<sup>8</sup>) trebuie să fie aleasă în valorile 89.55 pentru primul și, respectiv, 90.29 pentru cel de-al doilea eșantion de date. Pentru aceste valori ale pragului de decizie clasificatorul obținut este unul optimal.

Se observă că prin utilizarea unui set de date mult diminuat (doar 1% din primul eșantion de date) poziția suprafeței de decizie obținută pentru această aplicație se modifică, însă nu în mod substanțial (90.29 versus 89.55), eroarea fiind una acceptabilă.

- b) Dacă particularizăm problema și presupunem că  $\hat{\sigma}_1 = \hat{\sigma}_2 = \sigma$  atunci ecuația de gradul doi se transformă într-o ecuație de gradul întâi a cărei soluție este:

<sup>8</sup> Suprafețele de decizie sunt întotdeauna  $(d-1)$  dimensionale, vezi și **Subcapitolul 3.5.3**. Unde  $d$  este dimensiunea spațiului de intrare, a spațiului trăsăturilor.

$$a_d = \frac{\hat{m}_1 + \hat{m}_2}{2} + k \quad (7.26)$$

În relația (7.26),  $k$  depinde de raportul probabilităților apriorice ale claselor precum și de mediile acestora, fiind dat de:

$$k = \frac{\sigma^2}{\hat{m}_2 - \hat{m}_1} \ln \frac{P(c_1)}{P(c_2)} \quad (7.27)$$

- c) În momentul în care avem  $\hat{\sigma}_1 = \hat{\sigma}_2 = \sigma$  și, în plus,  $P(c_1) = P(c_2)$  – cu alte cuvinte, varianțele ambelor clase sunt aceleași iar clasele sunt echiprobabile (se particularizează relația (7.27) pentru  $P(c_1) = P(c_2)$ , rezultând  $k = 0$ ) – suprafața de decizie va fi poziționată la jumătatea distanței dintre centrele celor două clase. Altfel spus, suprafața de decizie este dată în această situație doar de medile claselor. Acest rezultat este identic cu cel ce se obține în cazul utilizării clasificatorului de minimă distanță bazat pe norma Euclidiană. În consecință, **în aceste condiții particulare clasificatorul Bayes-ian este identic cu un clasificator de tip minimă distanță**. Reamintim că modul de funcționare a clasificatorului de tip minimă distanță (vezi **Subcapitolul 4.2**) presupune utilizarea în procesul de clasificare numai a informației privind distanțele către mediile claselor.
- d) Pentru clasificatorul Bayes-ian, soluția ecuației suprafeței de decizie (vezi și relația (7.27) obținută pentru aplicația particulară, unidimensională, cu două clase, de mai sus), este una ce depinde, pe lângă informația dată de mediile claselor, și de informațiile furnizate de varianțele (mai general, matricile de covarianță) și, respectiv, probabilitățile apriorice ale claselor. Pentru a înțelege modul cum influențează aceste informații poziționarea suprafeței de decizie facem următoarele discuții pe cazul unidimensional al problemei noastre, extrapolarea acestor concluzii la cazul multidimensional fiind una directă:
- (i) Existența a două varianțe diferite pentru cele două clase ( $\hat{\sigma}_1 \neq \hat{\sigma}_2$ ) determină *mutarea pragului de decizie spre clasa cu varianță minimă* (respectiv, spre dreapta în cazul problemei în discuție; vezi **Figura 7.4**).

Acest fenomen este și unul intuitiv deoarece varianță minimă înseamnă, în principal, o concentrare mai mare a datelor în jurul valorii medii. În consecință, alegerea suprafeței de decizie

depinde și de varianța fiecărei aglomerări de puncte aparținând unei anumite clase nu numai de poziționarea centrului, a mediei clasei respective. Prin această observație evidențiem practic faptul că în procesul de clasificare *avem nevoie de o metrică care să depindă nu numai de distanța dintre elementul necunoscut la centrele claselor, ci și de varianța acestora* (de modul de distribuție al claselor). O astfel de metrică este, așa după cum știm, metrica Mahalanobis.

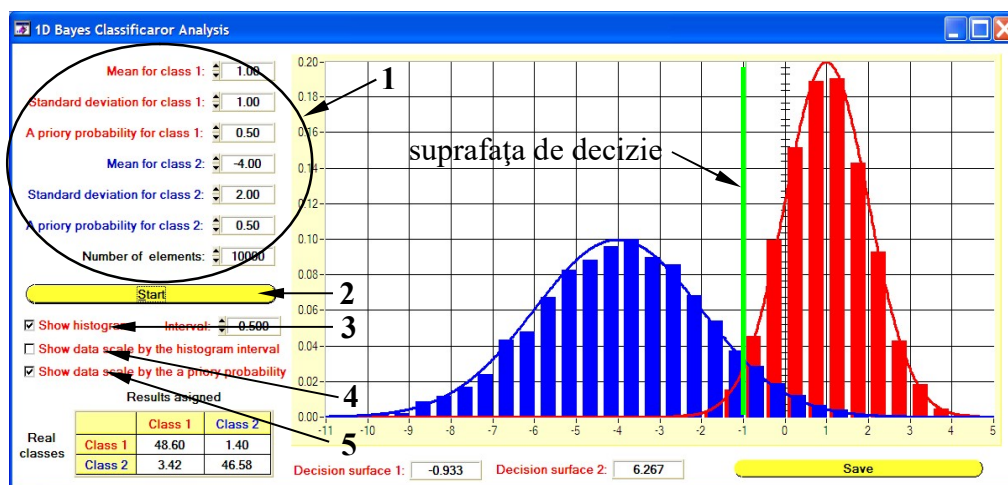
În contextul de mai sus, în cazul în care clasele sunt echiprobabile clasificatorul Bayes-ian este identic, din punct de vedere conceptual, cu un *clasificator de tipul Mahalanobis*.

- (ii) Din relația (7.27) se observă că raportul probabilităților apriorice ale claselor mută și el suprafața de decizie spre dreapta sau spre stânga, funcție de valoarea acestui raport. Intuitiv, ca regulă generală, pragul de decizie se mută întotdeauna mai aproape de clasa cu probabilitate mai mică. În consecință, clasa cu probabilitate mai mare va avea o regiunea de decizie mai mare.

Dacă luăm în discuție complexitatea diferiților clasificatori prezentați până acum observăm că:

- (a) clasificatorul de minimă distanță ce utilizează norma Euclidiană are complexitatea cea mai redusă, folosind în procesul decizional numai informația furnizată de *mediile claselor*;
- (b) clasificatorul bazat pe norma Mahalanobis folosește în plus în procesul de decizie, față de clasificatorul anterior, și informații ce țin de *dispunerea spațială a claselor*; aceste informații sunt valorificate prin intermediul matricii de covarianță (pentru spații de trăsături cel puțin bidimensionale), respectiv, prin intermediul varianței (pentru cazul unidimensional). Din acest punct de vedere complexitatea acestui clasificator este și ea una mai mare comparativ cu complexitatea clasificatorului de tip minimă distanță, studiat în Capitolul 4.
- (c) clasificatorul Bayes-ian utilizează în plus față de clasificatorul Mahalanobis – în determinarea funcțiilor discriminant și, deci, a poziției suprafețelor de decizie –, informații ce privesc *probabilitățile apriorice ale claselor*. În concluzie, acest din urmă clasificator utilizează în procesul decizional, simultan, informații ce privesc centrele claselor, dispunerea spațială a lor, precum și probabilitatea apriorică a acestora. Din acest punct de vedere, clasificatorul Bayes-ian este cel mai complex din cei prezentați aici și, mai mult, atât clasificatorul de minimă distanță cât și clasificatorul Mahalanobis reprezintă particularizări, în anumite

condiții bine definite (vezi discuțiile anterioare), ale clasificatorului *Bayes*-ian.



**Figura 7.6.** Interfața grafică a programului de analiză a clasificatorului *Bayes*-ian unidimensional

**Aplicație 7.2:** În directorul „*Bayes 1D*”, asociat acestui capitol, se găsește codul sursă al unei aplicații *software* dezvoltată în mediul LabWindows CVI capabilă:

1. Să genereze două seturi de date distribuite conform a două funcții densitate de probabilitate unidimensionale,  $f(x|C_1)$  și  $f(x|C_2)$ , ale căror parametri pot fi selectați de către utilizator (prin intermediul elementelor „1” de pe interfața grafică). Generarea unui nou set de date, cu noi parametri, se realizează prin apăsarea butonului „2”.
2. Să afișeze cele două funcții densitate de probabilitate:  $f(x|c_1)$  și  $f(x|c_2)$ .
3. Să afișeze histogramele celor două seturi de date prin bifarea controlului „3”.
4. Să afișeze cele două funcții densitate de probabilitate înmulțite cu intervalul pe care s-au determinat cele două histograme:  $f(x|C_1) \cdot \Delta x$  și  $f(x|C_2) \cdot \Delta x$  – similar cu reprezentarea grafică din **Figura 7.4**. Această afișare se realizează prin selectarea elementului de control „4”.
5. Să prezinte suprafața/suprafețele decizionale.
6. Să afișeze cele două funcții densitate de probabilitate înmulțite cu probabilitatea apriorică a celor două clase:  $f(x|c_1) \cdot P(c_1)$  și  $f(x|c_2) \cdot P(c_2)$ .

·  $P(c_2)$ . Acest mod de afișare este similar cu cel din **Figura 7.5** și se realizează prin selectarea elementului „5”.

7. Să prezinte funcțiile densitate de probabilitate scalate prin ambele informații prezentate la subpunctele 4 și 6 anterior prezentate.

Utilizând acest program se cere să se parcurgă următoarele cerințe:

1. Pentru două clase echiprobabile și având deviații standard egale, să se determine numărul soluțiilor obținute. De ce au fost obținute atâtea soluții? Unde este poziționată suprafața de decizie? De ce?
2. Păstrând parametri de la punctul 1 constanți, variați doar valoarea probabilității apriorice a uneia dintre clase. În ce mod se modifică poziția suprafeței de decizie? Către ce clasă? (*Notă:* după fiecare modificare a valorii probabilității apriorice a uneia dintre clase apăsați pe butonul „2” pentru vizualizarea rezultatelor).
3. Folosind aceleași valori ale parametrilor de la punctul 1, variați acum doar deviația standard a uneia sau alteia dintre clase și observați simultan modificarea poziției suprafeței/suprafețelor decizionale. În ce condiții/condiție clasificatorul *Bayes*-ian va genera două suprafețe de decizie?
4. Păstrând constante valorile parametrilor de la punctul 1, variați acum media uneia dintre clase. Observați variația performanțelor de clasificare obținute de clasificatorul *Bayes*-ian pentru valori mai mari sau mai mici ale mediei clasei modificate comparativ cu situația inițială.
5. Corelați modificarea suprafeței/suprafețelor decizionale generată de modificările de la subpunctele 2 și 3 cu performanțele de clasificare vizualizate cu ajutorul matricii confuziilor prezentă pe interfața grafică a programului.

### 7.2.2 Regula de decizie

Utilizând principiile generale, precum și relațiile prezentate anterior, ne este foarte ușor acum să construim un clasificator. Pentru a selecta clasa optimă putem alege una din următoarele două abordări:

- (a) Având un vector de trăsături  $a^0$ , calculăm mai întâi setul de probabilități posterioare:

$$P(c_i | a^0) = \frac{f_x(a^0 | c_i)P(c_i)}{P(a^0)} \quad (7.28)$$

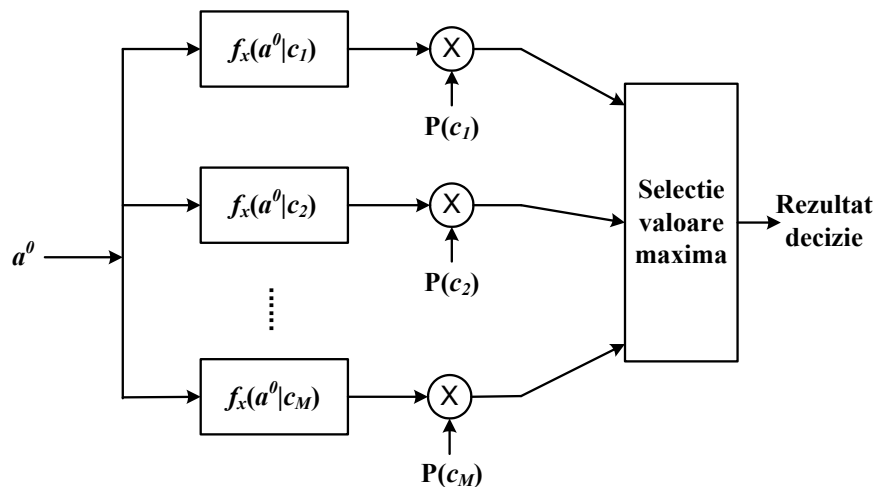
pentru toate clasele posibile și, ulterior, realizăm asignarea elementului  $a^0$  la acea clasă pentru care valoarea probabilității posterioare este maximă, conform regulei:

$$a^0 \text{ aparține clasei } c_i \text{ dacă } P(c_i | a^0) \geq P(c_j | a^0) \text{ pentru } \forall j \neq i \quad (7.29)$$

Schema bloc generală de implementare a clasificatorului *Bayesian* pentru  $M$  clase și un vector de intrare,  $a^0$ , multidimensional este prezentată în **Figura 7.7**.

- (b) În mod alternativ, putem determina mai întâi suprafața/suprafețele de decizie ce separă cele  $M$  clase și, ulterior, funcție de poziționarea vectorului de trăsături  $a^0$  față de aceste suprafețe de decizie, luăm decizia privind clasa optimă de apartenență a lui.

Astfel, spre exemplu, pentru **Problema 7.3**, se calculează mai întâi pragul  $a_d$  și ulterior, tragem concluzia de „subiect odihnit” ( $a^0 < a_d$ ) sau de „subiect obosit” ( $a^0 > a_d$ ).



**Figura 7.7.** Schema bloc a clasificatorului *Bayesian*

### 7.2.3 Eroarea de clasificare

În general, asignarea elementelor la una din cele două sau mai multe clase posibile nu se poate realiza chiar fără nici un fel de eroare.



Exemplificând pe **Problema 7.3** (vezi și **Figura 7.5**), se poate vedea cum “coada” distribuției de probabilitate pentru clasa indivizilor odihniți, ce se extinde și în dreapta punctului de decizie, ne dă eroarea de clasificare pentru clasa  $\{\text{odihnit}\}$  în timp ce “coada” distribuției de probabilitate pentru clasa indivizilor oboșiți, ce se extinde în stânga punctului de decizie, ne dă eroarea de clasificare pentru clasa  $\{\text{obosit}\}$ . În aceste condiții, *eroarea de clasificare globală* este dată de suma celor două erori, respectiv, de *suma suprafețelor mărginite de capetele funcțiilor de distribuție,  $f_x(a|c_i)$ , axă și pragul ales scalate la probabilitatea clasei respective*.

**Observație 7.7:** Cu cât suprapunerea este mai redusă cu atât eroarea de clasificare este și ea mai mică. Întrucât regiunile de decizie depind de pragul  $a_d$  ales, rezultă că și erorile de clasificare depind de poziționarea acestui prag.

Ținând cont de faptul că în majoritatea aplicațiilor practice, reale, eroarea de clasificare obținută cu diverși clasificatori propuși teoretic este una diferită de zero rezultă că rezolvarea unei probleme de clasificare revine la a găsi practic acel clasificator care să fie optimal.

Prin **clasificator optimal se înțelege acel clasificator pentru care probabilitatea de eroare de clasificare este minimă**.

Având în vedere faptul că modul de determinare a clasei dat de relația (7.13) minimizează probabilitatea de eroare, rezultă imediat de aici și faptul că ***clasificatorul Bayes-ian este clasificator optimal***.

**Observație 7.8:** Atenție! Reamintim aici că, prin clasificator optimal nu înțelegem un clasificator pentru care nu vom obține nici o eroare de clasificare ci vom înțelege un clasificator cu care obținem numărul minim de erori de clasificare posibile.

În cele ce urmează vom face o discuție privind eroarea de clasificare pentru cazul unidimensional, cu două clase posibile, și cu:

- (i) o singură soluție pentru ecuația suprafeței de decizie și, respectiv,
- (ii) două soluții pentru ecuația suprafeței de decizie.

**Discuții:**

- (i) Pentru a determina eroarea de clasificare pentru **Problema 7.3** – pentru care avem o singură soluție pentru ecuația suprafeței de decizie – va trebui să integrăm ariile,  $R_1$  și  $R_2$ , care ne dau tocmai această eroare (vezi **Figura 7.8**). Acest lucru este simplu într-un spațiu monodimensional dar devine dificil într-un spațiu multidimensional.

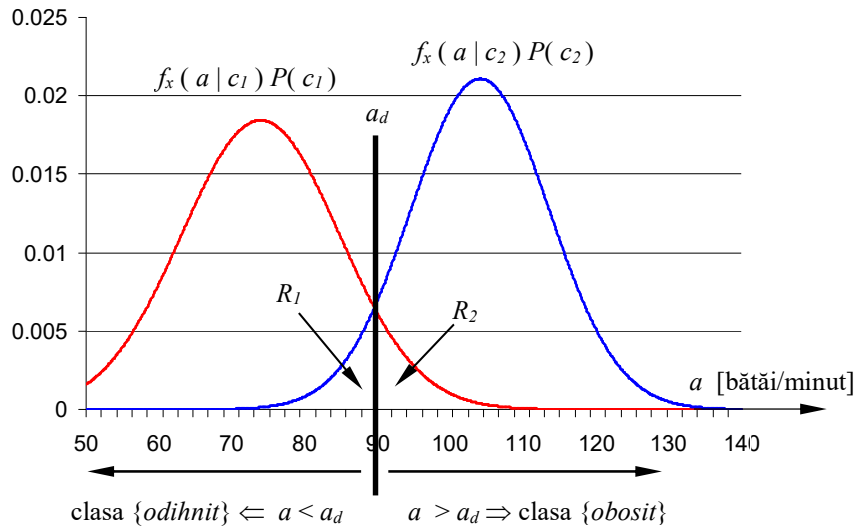
În cazul particular al problemei prezentate anterior probabilitatea de eroare este:

$$P_{eroare} = \int_{R_2} f_x(a | c_1)P(c_1)da + \int_{R_1} f_x(a | c_2)P(c_2)da \quad (7.30)$$

sau

$$P_{eroare} = \int_{a_d}^{+\infty} f_x(a | c_1)P(c_1)da + \int_{-\infty}^{a_d} f_x(a | c_2)P(c_2)da \quad (7.31)$$

Eroarea este cu atât mai mică (respectiv, acuratețea clasificatorului este mai mare) cu cât supunerea dintre clase este mai mică. În mod intuitiv, pentru varianțe egale ale claselor, cu cât va fi mai mare distanța dintre centrul clusterelor cu atât va fi mai mică suprapunerea între clase și, respectiv, cu atât va fi mai mică și eroarea de clasificare. În mod similar, pentru aceeași distanță între centrele claselor eroarea va fi cu atât mai mică cu cât varianța distribuției fiecărei clase, în parte, va fi mai mică.



**Figura 7.8.** Probabilitatea de eroare în cazul clasificatorului *Bayess-ian*

- (ii) În **Figura 7.9** este dată o reprezentare grafică pentru două clase arbitrare, caracterizate de două funcții densitate de probabilitate *gauss-iene*. În acest caz se poate observa că ecuația suprafeței de decizie are

două soluții,  $a_{d1}$  și  $a_{d2}$ . Pentru această situație eroarea de clasificare se calculează cu relația:

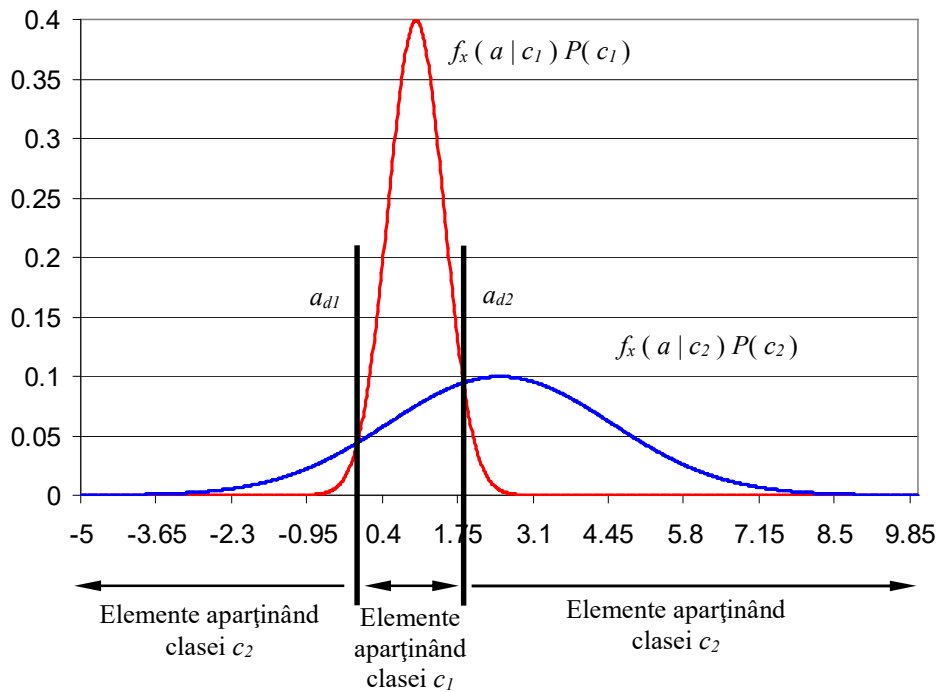
$$P_{eroare} = \int_{-\infty}^{a_{d1}} f_x(a | c_1)P(c_1)da + \int_{a_{d1}}^{a_{d2}} f_x(a | c_2)P(c_2)da + \int_{a_{d2}}^{+\infty} f_x(a | c_1)P(c_1)da \quad (7.32)$$

Mai general, probabilitatea medie a erorii de clasificare a clasificatorului Bayes-ian pentru cazul unei probleme cu două clase se scrie astfel:

$$P(eroare) = \int_{-\infty}^{+\infty} P(eroare, a)da = \int_{-\infty}^{+\infty} P(eroare | a)P(a)da \quad (7.33)$$

unde,

$$P(eroare | a) = \begin{cases} P(c_1 | a) & \text{dacă decidem clasa } c_2 \\ P(c_2 | a) & \text{dacă decidem clasa } c_1 \end{cases} \quad (7.34)$$



**Figura 7.9.** Probabilitatea de eroare în cazul clasificatorului Bayes-ian pentru două distribuții arbitrare și pentru două suprafețe de decizie

*În concluzie, factorii ce afectează eroarea de clasificare sunt, de fapt, o*

*combinație între diferența mediilor claselor, varianțele claselor și probabilitățile apriorice ale acestora.*

#### 7.2.4. Un exemplu bidimensional de recunoaștere de *pattern-uri*

Exemplul prezentat anterior (unidimensional) este prea simplu pentru a pune în evidență:

- întreaga metodologie utilizată în determinarea poziției suprafețelor de decizie – bazată pe modelarea statistică a setului de date;
- varietatea și caracteristicile suprafeței de decizie;
- anumite detalii și dificultăți ale proiectării clasificatorului.

Din aceste motive, în acest subcapitol vom trata cazul unui clasificator *Bayes-ian* bidimensional. Acest caz bidimensional a fost ales în principal datorită posibilității facile de vizualizare a seturilor de date și de înțelegere intuitivă a fenomenelor implicate în procesul clasificării. Desigur, această metodă poate fi ulterior generalizată pentru orice dimensiune arbitrară,  $d$ , a vectorului de trăsături.

În cele ce urmează considerăm aceeași problemă **Problema 7.3** (clasificarea unei populații de subiecți în două clase, {*odihnit*} versus {*obosit*}, funcție de activitatea cardiacă) doar că de această dată activitatea cardiacă a inimii va fi reprezentată prin doi parametri: **numărul de bătăi pe unitatea de timp** (minut) și, respectiv, **presiunea sistolică**. Utilizăm aceste variabile aleatoare pentru construirea unui vector de trăsături aleator bidimensional ce va fi prezentat în continuare la intrarea clasificatorului *Bayes-ian*.

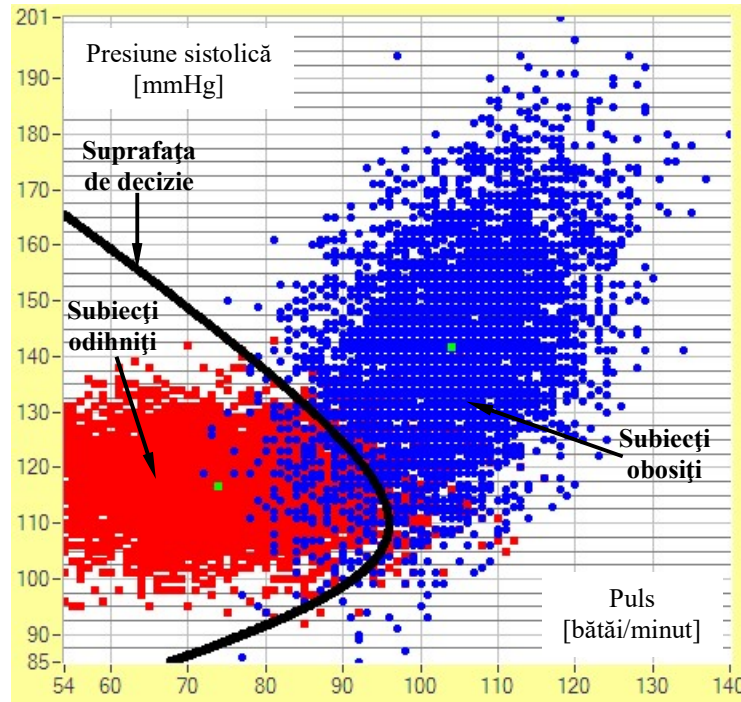
În continuare presupunem că funcțiile de distribuție ale celor două clase sunt *gauss-iene* iar probabilitățile apriorice ale claselor sunt ambele egale cu  $\frac{1}{2}$  (clasele sunt echiprobabile).

Scopul principal al acestui subcapitol este acela *de a determina poziția suprafeței de decizie astfel încât să obținem un clasificator optimal.*

**Notă:** Setul de date pe care îl vom prezenta și analiza în acest subcapitol se găsește în fișierul „**Odihnit-Obosit Puls-PresSist.txt**” din directorul asociat acestui subcapitol. O reprezentare grafică a unui eșantion de date de 10000 exemplare este redată, pentru această problemă, în **Figura 7.10**.

Conceptual metodologia de determinare a suprafețelor de decizie este similară ca și în cazul uni-dimensional al problemei anterior prezentate, unde suprafața de decizie s-a determinat din estimările făcute asupra mediei și deviației standard pentru cele două clase. În cazul multidimensional

diferențele sunt date de faptul că mediile claselor sunt vectori  $d$ -dimensionali iar varianțele claselor sunt matrici de covarianță,  $d \times d$  dimensionale.



**Figura 7.10.** Distribuția indivizilor funcție de activitatea cardiacă (puls și presiune sistolică) precum și suprafața optimală de decizie

În modul cel mai general suprafața de decizie dintre două clase  $i$  și  $j$  se găsește prin egalarea funcțiilor discriminant ce caracterizează clasele (vedeți **Subcapitolele 3.5.2 și 3.5.3**):

$$g_i(a) = g_j(a) \quad (7.35)$$

unde prin  $g$  s-au notat funcțiile discriminant ce caracterizează cele două clase.

În clasificatorul *Bayes*-ian regula de decizie este una exprimată, am văzut, în termenii probabilității posterioare a clasei (respectiv, probabilitatea clasei de a include un element particular analizat). Spre deosebitre de această formulare, o formă alternativă a regulii de decizie poate fi una exprimată în termenii funcțiilor discriminant, funcții ce atribuie un scor maxim elementelor clasei pe care o caracterizează și scoruri inferioare pentru vectorii de trăsături aparținând celorlalte clase.

Făcând echivalența celor două forme de exprimare a regulei de decizie, obținem:

$$g_i(a) = P(c_i | a) \quad \text{sau} \quad g_i(a) = f_x(a | c_i) \cdot P(c_i) \quad (7.36)$$

unde:  $f_x(a|c_i)$  este probabilitatea elementului  $a$  condiționată de clasa  $c_i$ .

**Observație 7.9:** În această ecuație am omis termenul  $f_x(a)$  deoarece el este un factor comun tuturor discriminanților și nu afectează, în esență, forma sau plasarea suprafeței de decizie; din acest motiv poate fi ignorat în definirea funcției discriminant.

Existența exponenților în definirea lui  $f_x(a|c_i)$  ne sugerează o alternativă de redefinire a funcțiilor discriminant prin intermediul logaritmului natural al relației (7.36). În acest context, noile relații vor fi date de:

$$g_i(a) = \ln f_x(a | c_i) + \ln P(c_i) \quad (7.37)$$

Relația de mai sus reprezintă forma cea mai generală pentru funcția discriminant, ea depinzând de funcția de distribuție a clasei  $c_i$  și de probabilitatea apriorică a clasei,  $P(c_i)$ .

Pentru vectorul aleator bidimensional real  $x$  ( $x = [x_1, x_2]^T$ ,  $x_1$  – fiind prima trăsătură – pulsul în cazul nostru particular –, în timp ce  $x_2$  este presiunea sistolică) funcția densitate de probabilitate pentru clasa  $i$  este:

$$f_x(a|c_i) = \frac{1}{(2\pi)^{D/2} |C_i|^{1/2}} e^{-\frac{1}{2}(a-m_i)^T C_i^{-1}(a-m_i)} \quad (7.38)$$

unde indicele  $i \in \{od, ob\}$ , indexează clasa  $od - \{odihnit\}$ , respectiv, clasa  $ob - \{obosit\}$ .

După logaritizarea relației (7.38), funcția discriminant devine:

$$g_i(a) = -\frac{1}{2}(a - m_i)^T C_i^{-1}(a - m_i) - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln|C_i| + \ln P(c_i) \quad (7.39)$$

Ținând cont de faptul că cele două clase sunt echiprobabile, ultimul termen din ecuația (7.39) va fi același pentru fiecare dintre cele două funcții discriminant în parte, el putând fi astfel eliminat. Același lucru este valabil și pentru termenul  $(D/2) \ln(2\pi)$ .

Din relația (7.39) se observă că, în mod similar problemei de clasificare uni-dimensionale, ceea ce contează în definirea funcției discriminant este distanța dintre element și media clasei, normalizată la matricea de covarianță. Analizând ecuațiile funcțiilor discriminant pentru fiecare clasă, în scopul găsirii suprafețelor de decizie, se observă că soluția căutată este, de fapt, o funcție dată de o distanță normalizată (distanță numită distanță

*Mahalanobis*<sup>9)</sup>, de matricile de covarianță și de probabilitățile apriori ale celor două clase.

**Tabelul 7.1.** Estimarea vectorilor medii și a matricilor de covarianță ale claselor pentru două seturi de date de dimensiuni diferite

	10000 de elemente	100 de elemente
<b>Clasa {odihnit}</b>	$m_{od} = \begin{bmatrix} 739 \\ 116.7 \end{bmatrix}$ $C_{od} = \begin{bmatrix} 117 & -4 \\ -4 & 53.2 \end{bmatrix}$	$m_{od} = \begin{bmatrix} 72.8 \\ 114 \end{bmatrix}$ $C_{od} = \begin{bmatrix} 114.4 & -12.6 \\ -12.6 & 53.9 \end{bmatrix}$
<b>Clasa {obosit}</b>	$m_{ob} = \begin{bmatrix} 104 \\ 141.7 \end{bmatrix}$ $C_{ob} = \begin{bmatrix} 89.5 & 70 \\ 70 & 277.4 \end{bmatrix}$	$m_{ob} = \begin{bmatrix} 102.8 \\ 142.8 \end{bmatrix}$ $C_{ob} = \begin{bmatrix} 108.9 & 86.3 \\ 86.3 & 300.3 \end{bmatrix}$

În tabelul de mai sus se prezintă valorile estimate pentru covarianțele și mediile claselor atunci când sunt luate în considerare 100, respectiv, 10000 de realizări particulare ale vectorului aleator pentru fiecare clasă în parte.

Suprafața de decizie a clasificatorului optimal (clasificatorul *Bayes*-ian) se obține în aceeași manieră ca și în cazul unidimensional, prin înlocuirea mediilor și a matricilor de covarianță estimate din datele de intrare – informații prezentate în **Tabelul 7.1** –, în relația (7.39) și egalarea funcțiilor discriminant pentru cele două clase.

Inversele matricilor de covarianță pentru ambele clase sunt următoarele:

$$C_{od}^{-1} = \begin{bmatrix} 8.5 \cdot 10^{-3} & 6.5 \cdot 10^{-4} \\ 6.5 \cdot 10^{-4} & 1.8 \cdot 10^{-2} \end{bmatrix} \quad (7.41)$$

<sup>9)</sup> Reamintim că această distanță normalizată este dată de:

$$d^2 = (a - \hat{m}_x)^T \hat{C}_x^{-1} (a - \hat{m}_x) \quad (7.40)$$

În relația (7.40)  $\hat{m}_x$  reprezintă vectorul mediu estimat al clasei iar  $\hat{C}_x$  reprezintă estimatul matricei de covarianță.

$$C_{ob}^{-1} = \begin{bmatrix} 1.4 \cdot 10^{-2} & -3.5 \cdot 10^{-3} \\ -3.5 \cdot 10^{-3} & 4.5 \cdot 10^{-3} \end{bmatrix} \quad (7.42)$$

Determinanții acestor matrici sunt următorii  $\det(C_{od}) = 6214.7$  și  $\det(C_{ob}) = 19926$ . Introducând toate aceste valori în (7.39) obținem:

$$g_{od}(a) = -\frac{1}{2}[a_1 - 73.9 \cdot a_2 - 116.7] \cdot \begin{bmatrix} 8.5 \cdot 10^{-3} & 6.5 \cdot 10^{-4} \\ 6.5 \cdot 10^{-4} & 1.8 \cdot 10^{-2} \end{bmatrix} \cdot \begin{bmatrix} a_1 - 73.9 \\ a_2 - 116.7 \end{bmatrix} - \frac{1}{2} \ln(6214.7) \quad (7.43)$$

Pentru clasa subiecților oboșiți discriminantul este:

$$g_{ob}(a) = -\frac{1}{2}[a_1 - 104 \cdot a_2 - 141.7] \cdot \begin{bmatrix} 1.4 \cdot 10^{-2} & -3.5 \cdot 10^{-3} \\ -3.5 \cdot 10^{-3} & 4.5 \cdot 10^{-3} \end{bmatrix} \cdot \begin{bmatrix} a_1 - 104 \\ a_2 - 141.7 \end{bmatrix} - \frac{1}{2} \ln(19926) \quad (7.44)$$

Suprafața de decizie se obține egalând funcțiile discriminant ale celor două clase:

$$g_{od}(a) = g_{ob}(a) \quad (7.45)$$

Ecuția suprafeței de decizie ce se obține în urma calculelor este :

$$a_1^2 - 8.15 a_2^2 + 25228.58 a_1 + 20895.12 a_2 + 177.6 a_1 a_2 - 2808973.26 = 0 \quad (7.46)$$

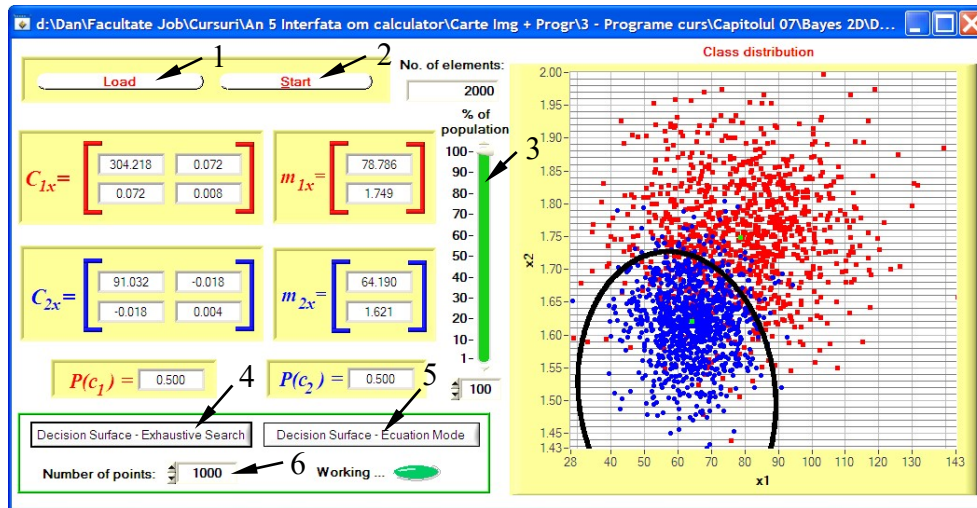
Această suprafață este o quadratică în spațiul bidimensional al trăsăturilor, suprafață ce determină obținerea celei mai mici erori de clasificare pentru problema de clasificare abordată aici. O reprezentare grafică a celor două clase precum și a suprafeței de decizie optimale obținute este dată în **Figura 7.10**.

**Problemă 7.4:** Pornind de la relația (7.39) să se determine, pentru cazul bidimensional al vectorului aleator de trăsături  $x$ , formula generală ce descrie ecuația suprafeței de decizie (formulă similară relației (7.46)). Implementați această relație într-un program scris în limbajul LabWindows CVI și verificați corectitudinea ei prin reprezentarea grafică a suprafeței de decizie.

Din simpla vizualizare a claselor și, respectiv, a poziției suprafeței de decizie, vezi **Figura 7.10**, observăm faptul că, chiar și în aceste condiții, a utilizării unui clasificator optimal, se pot obține, așa cum am menționat deja anterior, foarte multe erori de clasificare. Reamintim aici faptul că **un clasificator optimal**, ai cărui parametri sunt corect estimați, **va obține – dintre toate familiile de clasificatori existenți –, nu eroare zero ci cea mai mică eroare de clasificare pentru problema de clasificare abordată.**



În spații multidimensionale, estimarea corectă, cu o precizie adecvată, a tuturor elementelor matricelor de covarianță aferente claselor este, în general, o sarcină dificilă. În practică, matricele de covarianță, precum și mediile claselor, sunt adesea estimate nesatisfăcător, aceasta și datorită unor seturi de date insuficient de mari, ceea ce conduce, în mod direct, și la obținerea unor suprafețe de decizie cu performanțe suboptimale.



**Figura 7.11.** Interfața grafică a programului utilizat în clasificarea unor seturi de date caracterizate de vectori bidimensionali de trăsături.

**Aplicație 7.3:** În cadrul acestei exemplu se analizează, cu ajutorul unei aplicații practice, un clasificator *Bayes*-ian bidimensional. Codul programului și *kit*-ul de instalare se găsesc în directorul „*Bayes 2D*” asociat acestui capitol. Acest program permite încărcarea oricărui set de date organizat pe trei coloane. Pe primele două coloane sunt valorile numerice ale celor două trăsături ce formează vectorul de trăsături, în timp ce elementele de pe ultima coloană pot lua doar valori egale cu 0 sau 1 și furnizează informația de apartenență a elementului la una din cele două clase. Încărcarea setului de date se realizează prin apăsarea butonului „1”, vezi **Figura 7.11**. O dată cu încărcarea setului de date acesta se va afișa automat. Prin apăsarea butonului „2” se inițiază procesul de calcul a parametrilor clasificatorului *Bayes*-ian, parametri ce sunt ulterior prezentați pe interfața grafică. Cu ajutorul acestui program se cere atingerea următoarelor obiective:

1. Încărcați unul din următoarele seturi de date: “*Odihnit-Obosit Puls-PresSist rnd.txt*”, “*Odihnit-Obosit Puls-PresSist.txt*”, “*Femei-Barb Greut-Inalt Bayes rnd.txt*”, “*Femei-Barb Greut-Inalt Bayes.txt*” sau “*Femei-Barb Inalt-Greut Bayes.txt*”.
2. Apăsați butonul „Start” pentru determinarea parametrilor statistici ai celor două clase. Verificați în mod intuitiv corelația existentă între dispunerea spațială a celor două clase și coeficienții matricilor de covarianță obținute.
3. Verificați influența mărimii setului de date asupra estimării matricilor de covarianță și a vectorilor medii. Puteți modifica mărimea setului de date utilizat de clasificator prin schimbarea valorii controlului „3” de pe interfața grafică. Este recomandat ca în cadrul acestui subpunct să utilizați unul din seturile de date ce au vectorii de trăsături aparținând celor două clase amestecați într-un mod aleator (fișierele ce conțin aceste seturi de date au incluse în numele lor caracterele „*rnd*”).
4. Reprezentați grafic suprafața de decizie dintre cele două clase. În cadrul programului sunt implementate două metode de determinare a acestei suprafețe de decizie. În prima din ele, se realizează o parcurgere exhaustivă a spațiului de intrare într-un număr de pași selectabili din controlul „6”. În cea de a doua metodă se determină ecuația suprafeței de decizie, obținându-se o relație similară cu (7.46), care ulterior este reprezentată. Înțelegeți codurile acestor două subrutine.
5. Observați variabilitatea acestor suprafețe de decizie funcție de mărimea setului de date.

**Aplicație 7.4:** Pentru programul prezentat în cadrul **Aplicației 7.1** parcurgeți pașii 1-5 încă o dată însă, de această dată, folosind și clasificatorul *Bayes*-ian.

Utilizând setul de date „*Femei-Barb Greut-Inalt Bayes rnd.txt*” trasați suprafețele de decizie generate de clasificatorii *Mahalanobis* și *Bayes*-ian și analizați, folosind matricea confuziilor, performanțele obținute.

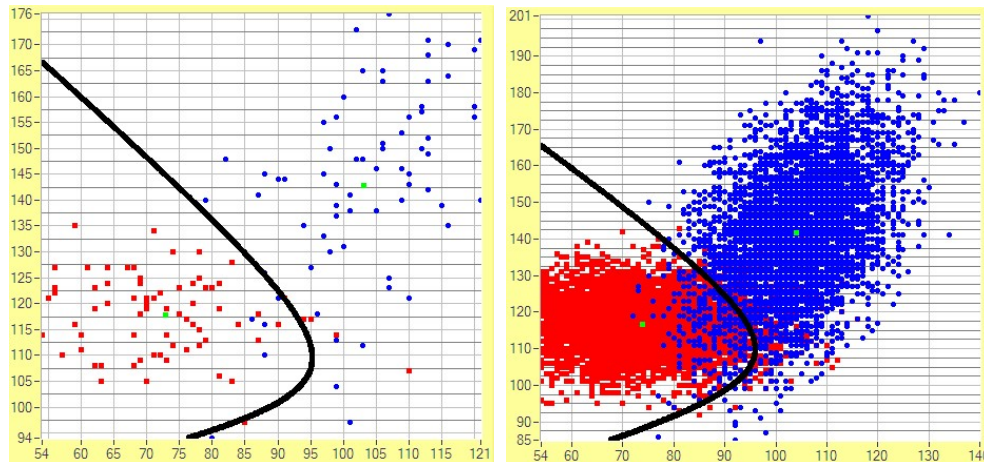
Trasați simultan suprafețele de decizie pentru clasificatorul *Mahalanobis* și pentru clasificatorul *Bayes*-ian utilizând, într-o primă fază setul de date „*6 Classes.txt*” iar, ulterior, setul de date „*6 Classes - Different prob.txt*”. Puneți în evidență modificările ce au loc în

cadrul clasificatorului *Bayes*-ian atunci când probabilitățile apriorice ale claselor nu mai sunt egale.

### 7.2.5. Sensibilitatea funcțiilor discriminant funcție de mărimea setului de date

În subcapitolele precedente am prezentat o metodă cu ajutorul căreia am reușit să determinăm suprafața optimă de decizie în ipoteza că funcția de distribuție de probabilitate pentru fiecare clasă este de tip densitate *gauss*-iană. Această metodă de clasificare este una foarte puternică însă ea are și următoarele dezavantaje:

1. se bazează pe anumite ipoteze în legătura cu distribuția setului de date, ipoteze asupra formei funcțiilor densitate de probabilitate ce caracterizează setul de date și, în plus,
2. metoda necesită un volum mare de date pentru estimarea, cu erori mici, a parametrilor funcțiilor discriminant.



**Figura 7.12.** Poziționarea suprafețelor de decizie pentru mărimi diferite ale seturilor de date (**a**). 100 de elemente, și (**b**). 10000 de elemente.

În consecință, și *calitatea rezultatelor obținute depinde*, pe de o parte, de *validitatea ipotezelor făcute*, iar pe de altă parte, de *calitatea valorilor estimate obținute pentru parametrii populației* – calitate ce depinde direct de mărimea eșantionului de date.

În acest subcapitol vom ilustra efectul mărimii setului de date asupra corectitudinii estimării parametrilor funcțiilor discriminant.

**Exemplul 7.3:** Să presupunem că avem numai 100 de realizări particulare (vectori bidimensionali de trăsături:  $a = [\text{puls}, \text{presiune sistolică}]^T$ ) din care 50 de elemente corespund unor subiecți aleși în mod aleator din clasa  $\{\text{odihnit}\}$  iar 50 de vectori de trăsături aparțin unor subiecți aleși în mod aleator din clasa  $\{\text{obosit}\}$ . Dacă calculăm media și matricea de covarianță pentru fiecare clasă în parte obținem rezultatele din **Tabelul 7.1**. Din analiza acestui tabel se observă că valorile diferiților parametri se modifică față de cele calculate pe setul de 10000 de vectori de trăsături. În momentul când construim, pe baza acestor parametri, funcția discriminant optimă pentru eșantionul de 100 de elemente observăm că forma ei este una similară cu cea obținută pe eșantionul mai mare de date, vezi **Figura 7.12**, însă poziția în spațiul de intrare se modifică, ceea ce duce în mod corespunzător și la modificarea probabilității de eroare (în particular, avem de-a face, așa cum era de așteptat, cu o creștere a acestei probabilități).

**Tabelul 7.2.** Estimarea parametrilor statistici pentru diferite mărimi ale seturilor de date

	100 de subiecți	1000 de subiecți
<b>Femei</b>	$m_f = \begin{bmatrix} 1.626 \\ 61.969 \end{bmatrix}$ $C_f = \begin{bmatrix} 0.004 & -0.182 \\ -0.182 & 104.9 \end{bmatrix}$	$m_f = \begin{bmatrix} 1.619 \\ 64339 \end{bmatrix}$ $C_f = \begin{bmatrix} 0.003 & -0.023 \\ -0.023 & 93 \end{bmatrix}$
<b>Bărbați</b>	$m_{ob} = \begin{bmatrix} 1.731 \\ 77.469 \end{bmatrix}$ $C_{ob} = \begin{bmatrix} 0.009 & 0.117 \\ 0.117 & 249.795 \end{bmatrix}$	$m_{ob} = \begin{bmatrix} 1.744 \\ 78.64 \end{bmatrix}$ $C_{ob} = \begin{bmatrix} 0.007 & 0.152 \\ 0.152 & 313.256 \end{bmatrix}$

Clasa cea mai afectată în sensul modificării valorilor estimate ale parametrilor este clasa subiecților oboșiți (modificări cu aproape cca 20 de puncte ale elementelor de pe diagonala principală a matricii de covarianță, atunci când se trece de la o analiză a unui eșantion format din 10000 de elemente la unul format din doar 100 de elemente).

**Exemplul 7.4:** Într-un alt exemplu (problema discriminării în două clase,  $\{\text{femei}\}$  și  $\{\text{bărbați}\}$ , a unui eșantion de subiecți, caracterizat de vectorul aleator  $x = [\text{înălțime}, \text{greutate}]^T$ ), se observă, conform **Tabelul 7.2**, că în cazul trecerii de la un eșantion inițial de 1000 de

subiecți (net mai mic față de cel prezentată în exemplul anterior, de 10000 de elemente) la unul de 100 de subiecți, se obțin, de exemplu, variații de cca 60 de unități în determinarea elementelor matricii de corelație pentru clasa {bărbați}.

***O primă concluzie ce s-ar putea desprinde din aceste două exemple ar fi aceea că „sensibilitatea” parametrilor clasificatorului Bayes-ian este puternic dependentă de tipul problemei analizate*** (de exemplu, de particularitățile trăsăturilor analizate) precum și de specificitățile setului de date.

O analiză suplimentară a modului cum se modifică probabilitatea de eroare la o modificare a dimensiunii eșantionului ne aduce un surplus informațional important și anume: scăderea dimensiunii eșantionului se însoțește de o creștere a erorii de clasificare și invers; o creștere a performanțelor clasificatorului atunci când estimarea parametrilor se face pe un eșantion mai mare de date se explică prin aceea că valoarea unui estimator, în general, tinde spre valoarea reală a parametrului populației atunci când dimensiunea eșantionului tinde, la limită, spre infinit.

Problema de clasificare abordată în ***Subcapitolul 7.2.4*** a fost una relativ foarte simplă, în doar două dimensiuni, 6 parametri trebuind să fie estimați pentru fiecare clasă (3 parametri ai matricii de covarianță, 2 parametri ai vectorului mediu și probabilitatea apriorică a clasei) din cele 50 de elemente din cât este formată o clasă. Într-un spațiu multidimensional se poate întâmpla ca numărul parametrilor ce trebuie să fie estimați să fie de același ordin de cardinalitate cu mărimea setului de date – în această situație suprafața de decizie va fi poziționată cu certitudine departe de cea optimă datorită, în principal, estimării nesatisfăcătoare a parametrilor clasificatorului. Pentru o analiză și descriere completă a acestei probleme citiți și subcapitolul în care se analizează avantajele și dezavantajele clasificatorului Mahalanobis (vezi ***Subcapitolul 7.1.2***).

***Găsirea, analiza și implementarea unor clasificatori care să fie mai puțin sensibili la estimarea parametrilor necunoscuți ai populației*** este un obiectiv principal a teoriei clasificatorilor.

În literatură s-a ajuns la concluzia că pentru a atinge acest obiectiv trebuie ca ***forma funcțională a funcțiilor discriminant să fie cât mai simplă*** [Principe, 2000]. ***În principal este de dorit să utilizăm funcții discriminant care au cât mai puțini parametri*** iar acești parametri să poată fi estimați într-un mod cât mai robust posibil din setul de date pe care îl avem la dispoziție. Aparent, astfel de funcții discriminant, mai simple, pot fi unele suboptimale pentru problema dată; experiența însă ne arată că, frecvent, cu aceste funcții suboptimale putem obține performanțe mai bune

decât cele generate de un clasificator, teoretic, optimal. Acest lucru poate părea un paradox dar explicația este una ce ține de estimarea, de cele mai multe ori, inexactă a parametrilor clasificatorului optimal. Astfel, chiar dacă folosim funcții discriminant cuadratică (considerate optimale pentru clase cu distribuție de probabilitate *gauss-iană*), similare cu cele din relația (7.46), acestea pot fi inexact poziționate, vezi **Figura 7.12(b)**, obținându-se astfel erori destul de mari de clasificare.

### 7.2.6. Selecția trăsăturilor bazată pe funcția densitate de probabilitate posterioară

Dintr-un proces oarecare se pot extrage un număr foarte mare de trăsături. O dată cu creșterea numărului de trăsături folosite, crește însă și complexitatea clasificatorului. De aceea se pune problema selecției acelor trăsături care sunt purtătoare de informație discriminantă maximă, necesară în procesul de clasificare.

Pentru un clasificator *Bayes-ian* se definește *riscul condiționat* ca fiind dat de relația:

$$R(c_i | a) = 1 - f_x(a | c_i)P(c_i) \quad (7.47)$$

La relația de mai sus s-a ajuns astfel:

- fie  $\{\alpha_1, \dots, \alpha_M\}$  mulțimea finită a deciziilor posibile, cu  $\{\alpha_1, \dots, \alpha_M\} \subseteq \{c_1, \dots, c_M\}$ ;
- se definește o *funcție de pierdere*,  $\lambda(\alpha_i | c_j)$ , ce stabilește *pierdere* pe care o atrage după sine o decizie greșită a unei clase  $\alpha_i$  în locul clasei  $c_j$  reale;
- presupunem că observăm o realizare particulară,  $a$ , a vectorului aleator multidimensional,  $x$ , pentru care luăm decizia de apartenență la clasa  $\alpha_i$ ; dacă adevărata clasă de apartenență a lui  $a$  este  $c_j$  atunci pierdere pe care o obținem în acest caz este  $\lambda(\alpha_i | c_j)$ ;
- întrucât  $P(c_j | a)$  este probabilitatea clasei adevărate atunci, pierdere medie (numită și *risc condiționat*, în teoria deciziei), generată de alegerea clasei  $\alpha_i$ , este dată de relația:

$$R(\alpha_i | a) = \sum_{j=1}^M \lambda(\alpha_i | c_j)P(c_j | a) \quad (7.48)$$

- pentru cazul particular al funcției de pierdere de tip zero-unu, definită astfel:

$$\lambda(\alpha_i | c_j) = \begin{cases} 0 & , \text{pentru } i = j \\ 1 & , \text{pentru } i \neq j \end{cases}, \text{ cu } i, j = \overline{1, M} \quad (7.49)$$

*riscul condiționat* definit prin relația (B041) devine chiar probabilitatea medie a erorii fiind dat de:

$$R(\alpha_i | a) = \sum_{j=1}^M \lambda(\alpha_i | c_j) P(c_j | a) = \sum_{j \neq i} P(c_j | a) = 1 - P(c_i | a) \quad (7.50)$$

în timp ce eroare globală de clasificare este dată de:

$$\varepsilon = \sum_{i=1}^c \sum_{j=1, j \neq i}^c \int (1 - R(c_i | a)) da \quad (7.51)$$

unde  $a$  este un vector de trăsături. Din relația de mai sus se poate scrie:

$$\varepsilon = \sum_{i=1}^c \sum_{j=1, j \neq i}^c P(c_i) \int f_x(a | c_i) da \quad (7.52)$$

O expresie asemănătoare se poate scrie și pentru fiecare trăsătură în parte:

$$\varepsilon_k = \sum_{i=1}^c \sum_{j=1, j \neq i}^c P(c_i) \int f_x(a_k | c_i) da_k \quad (7.53)$$

Pentru a selecta acele trăsături purtătoare de o cantitate maximă de informație discriminatorie se pot urma pașii: (1). pentru fiecare trăsătură în parte se calculează  $\varepsilon_k$ , cunoscând aprioric sau estimând  $P(c_i)$  și  $f_x(a_k | c_i)$ , iar apoi, (2). se rețin acele trăsături pentru care  $\varepsilon_k$  este minim (mai mic decât un anumit prag) și se elimină celelalte trăsături, considerându-se că sunt improprii pentru o clasificare corectă a setului de date.

### 7.3. Probleme

1.
  - (a) Care este ideea de bază, fundamentală, care stă în spatele clasificatorului *Bayes*-ian (exprimați această idee prin cuvinte)? Ce reprezintă fiecare termen din relația lui *Bayes* – relație utilizată în cadrul clasificatorului *Bayes*-ian?
  - (b) Clasificatorul *Bayes*-ian este considerat un clasificator optimal. Ce înțelegeți prin clasificator optimal?
  - (c) Explicați influența varianței trăsăturilor asupra poziționării suprafeței de decizie.
  - (d) Explicați influența probabilității apriorice a clasei asupra poziționării suprafeței de decizie.
  - (e) Explicați influența setului de date asupra performanțelor clasificatorului *Bayes*-ian.
  - (f) Desenați schema bloc a clasificatorului *Bayes*-ian.
  
2. Fie două clase distincte caracterizate de două funcții de densitate de probabilitate *gauss*-iene unidimensionale. Prima clasă este de medie 1.5 și varianță 0.04 în timp ce cea de a doua clasă are media 2 și varianța 0.64. Probabilitatea apriorică a primei clase este 1/5 iar a celei de a doua clasă 4/5.
  - (a) Desenați în același grafic ambele funcții densitate de probabilitate. Reprezentarea grafică va fi una calitativă dar va ține cont în reprezentare de parametrii celor două funcții de densitate.
  - (b) Scrieți relația matematică ce caracterizează funcția discriminant pentru elementele primei clase.
  - (c) Desenați schema bloc a clasificatorului *Bayes*-ian ce utilizează două funcții discriminant caracteristice fiecărei clase în parte.
  - (d) Determinați, în mod numeric, poziția exactă a suprafețelor de decizie.
  - (e) Realizarea particulară 0.5, a variabilei aleatoare  $x$ , aparține primei clase sau a celei de a doua ? Dar 1.4 ? Dar 1.8 ? Pentru fiecare valoare justificați, în mod numeric, decizia luată.
  - (f) Dacă cea de a doua clasă va avea o probabilitate apriorică mai mică ce se va întâmpla cu fiecare suprafață de decizie în parte? Justificați-vă răspunsul în mod analitic sau printr-o analiză conceptuală.



3. Să presupunem că un vector de trăsături mono-dimensional  $x$  este utilizat pentru a decide între clasa  $\omega_1$  și  $\omega_2$ . Probabilitatea a priori a clasei  $\omega_1$  este  $2/7$ . Funcțiile densitate de probabilitate pentru  $x$  în ipoteza apartenenței la clasele  $\omega_1$  și  $\omega_2$  sunt date de relațiile:

$$f_x(a|\omega_1) = \begin{cases} \frac{1}{4}, & -3 \leq a \leq 1 \\ 0, & \text{în rest} \end{cases} \quad (7.54)$$

$$f_x(a|\omega_2) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(a-m_x)^2}{2\sigma_x^2}} \quad (7.55)$$

- Demonstrați ca (7.54) este o funcție de tip densitate de probabilitate legitimă (verificați proprietățile acesteia).
  - Care este probabilitatea a priori pentru clasa  $\omega_2$ . Sunt clasele  $\omega_1$  și  $\omega_2$  echiprobabile?
  - Determinați pragul selecției optime a clasei  $\omega_2$  versus  $\omega_1$  (pentru situația în care  $\sigma_x = 1/2$  iar  $m_x = 2$ ).
  - Reprezentați grafic modalitatea de obținere a suprafeței de decizie în cazul utilizării unui clasificator *Bayes*-ian. Cum ar arăta această reprezentare grafică dacă  $m_x = 1.5$ ?
4. O variabilă aleatoare  $x$  este utilizată pentru a discrimina două clase,  $\omega_1$  și  $\omega_2$ . Pentru prima clasă știm că  $P(\omega_1) = 1/6$  iar pentru cea de a doua clasă  $P(\omega_2) = 5/6$ . Funcțiile densitate de probabilitate pentru variabila aleatoare  $x$  în ipoteza apartenenței la clasele  $\omega_1$  și  $\omega_2$  sunt date de:

$$f_x(a|\omega_1) = \begin{cases} \frac{1}{2}, & 2 \leq a \leq 4 \\ 0, & \text{în rest} \end{cases} \quad (7.56)$$

$$f_x(a|\omega_2) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(a-m_x)^2}{2\sigma_x^2}} \quad (7.57)$$

- Demonstrați că relația (1) este o funcție densitate de probabilitate legitimă (verificați proprietățile acesteia).
- Care este semnificația termenului echiprobabil? Sunt clase  $\omega_1$  și  $\omega_2$  echiprobabile?
- Determinați pragul selecției optime a clasei  $\omega_1$  versus  $\omega_2$  pentru situația în care  $\sigma_x = 1/2$  iar  $m_x = 1.5$ .

- (d). Reprezentați grafic funcțiile densitate de probabilitate ale celor două distribuții. Această reprezentare grafică va fi corelată cu rezultatele obținute la punctele anterioare.

Valori ce pot fi utilizate în determinările numerice:

$$\begin{array}{ll} \ln 1 = 0; & \ln 2 \approx 0.7; \\ \ln 3 \approx 1; & \ln 4 \approx 1.4; \\ \ln 5 \approx 1.6; & \ln 6 \approx 1.8; \\ \ln 7 \approx 1.9; & \ln 8 \approx 2; \\ \ln 9 \approx 2.2; & \ln 10 \approx 2.3. \\ \sqrt{2\pi} \approx 2.5 & \end{array}$$

5. O variabilă aleatoare  $x$  este utilizată pentru a discrimina două clase echiprobabile,  $\omega_1$  și  $\omega_2$ . Funcțiile densitate de probabilitate pentru variabila aleatoare  $x$  în ipoteza apartenenței la clasele  $\omega_1$  și  $\omega_2$  sunt date de:

$$f_x(a|\omega_1) = \begin{cases} \frac{1}{2}, & 0 \leq a \leq 2 \\ 0, & \text{în rest} \end{cases} \quad (7.58)$$

$$f_x(a|\omega_2) = \begin{cases} 0 & a < 1 \\ 0.4 \cdot a - 0.4 & 1 \leq a \leq 3 \\ 5.6 - 1.6 \cdot a & 3 < a \leq 3.5 \\ 0 & a > 3.5 \end{cases} \quad (7.59)$$

- (a) Demonstrați că relația (7.58) este o funcție densitate de probabilitate legitimă (verificați proprietățile acesteia).
- (b) Determinați pragul selecției optimale a clasei  $\omega_1$  versus  $\omega_2$ .
- (c) Generați o reprezentare grafică în care să se pună în evidență existența pragului de selecție ce separă cele 2 clase. Notați pe această prezentare grafică funcțiile pe care le reprezentați. Această reprezentare grafică va fi corelată cu rezultatul obținut la punctul anterior.
- (d) Determinați probabilitatea de eroare a acestui clasificator.
- (e) Determinați probabilitatea ca o valoare din eșantion să fie corect clasificată.
- (f) Demonstrați că relația (7.59) este o funcție densitate de probabilitate legitimă (verificați proprietățile acesteia).
- (g) Explicați în mod intuitiv ce se va întâmpla cu poziția suprafeței de decizie și de ce se va întâmpla acest lucru (comparativ cu situația

anterioară existentă în rezolvarea punctelor (a)-(f)) dacă probabilitățile apriorice ale claselor vor fi următoarele:  $P(\omega_1) = 2/6$  și, respectiv,  $P(\omega_2) = 4/6$ .

6. Două funcții densitate de probabilitate date de:

$$f_x(a|\omega_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{a^2}{2}} \quad (7.60)$$

$$f_x(a|\omega_2) = \frac{4}{\sqrt{2\pi}} e^{-(2a-1)^2} \quad (7.61)$$

caracterizează două clase ( $\omega_1$  și  $\omega_2$ ) ce au probabilitățile apriorice ale claselor  $P(\omega_1) = \frac{4}{5}$  și  $P(\omega_2) = \frac{1}{5}$ .

- Determinați poziția optimă a suprafeței de decizie folosindu-se pentru aceasta un clasificator *Bayes*-ian.
- Determinați probabilitatea de eroare a clasificatorului *Bayes*-ian pentru datele de la punctul anterior.
- Presupunând de această dată cele două clase ( $\omega_1$  și  $\omega_2$ ) echiprobabile, definite de următoarele funcții densitate de probabilitate:

$$f_x(a|\omega_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(a-4)^2}{2}} \quad (7.62)$$

și, respectiv,

$$f_x(a|\omega_2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(a-6)^2}{2}} \quad (7.63)$$

determinați poziția suprafeței de decizie optimale fără efectuarea nici unui calcul prealabil. Argumentați-vă decizia.

- Demonstrați invarianța distanței Mahalanobis  $r^2 = (x - m_x)^T C_x^{-1} (x - m_x)$  la orice transformare liniară de forma  $\mathbf{y} = \mathbf{A} \mathbf{x}$ .
  - Cum utilizați distanța Mahalanobis în cadrul unui clasificator de tip minimă distanță ?
  - Care sunt avantajele utilizării metricii Mahalanobis ?
  - Care este dezavantajul major al utilizării metricii Mahalanobis ? Discutați această problemă în paralel și în cazul clasificatorului *Bayes*-ian.
  - Care este ideea de bază în cazul clasificatorului *Bayes*-ian ?

- (f) În ce condiții un element va fi asignat unei clase în cazul clasificatorului *Bayes-ian* ? (prezentați inclusiv relația care asigură această asignare și explicați fiecare termen al acesteia).

