

# Diskriminantna analiza

Prof.dr.sc. Bojana Dalbelo Bašić  
Fakultet elektrotehnike i računarstva


## Uvod

Primjer:

Kako odrediti sigurnu vrstu anestetika za pacijenta koji mora na operaciju ?

O pacijentu su poznati određeni podaci, npr.: dob, spol, krvni tlak, težina, krvna slika, alergijska reakcija, itd.

Temeljem ovakvog znanja anesteziolog mora odrediti da li i koji anestetik smije dati pacijentu.



## Temeljem ovakvog znanja anesteziolog želi znati:

- Može li se oblikovati pravilo koje bi nove pacijente ispravno klasificiralo u dvije grupe: grupu onih kojima se smije dati anestetik, te grupu onih kojima se ne smije dati anestetik ?
- Ako se može oblikovati pravilo, koje je to pravilo ?
- Kolika je vjerojatnost krive klasifikacije?
- Koje su posljedice krive klasifikacije?
- Što najviše utječe na klasifikaciju?

Otkrivanje znanja u skupovima  
podataka

3



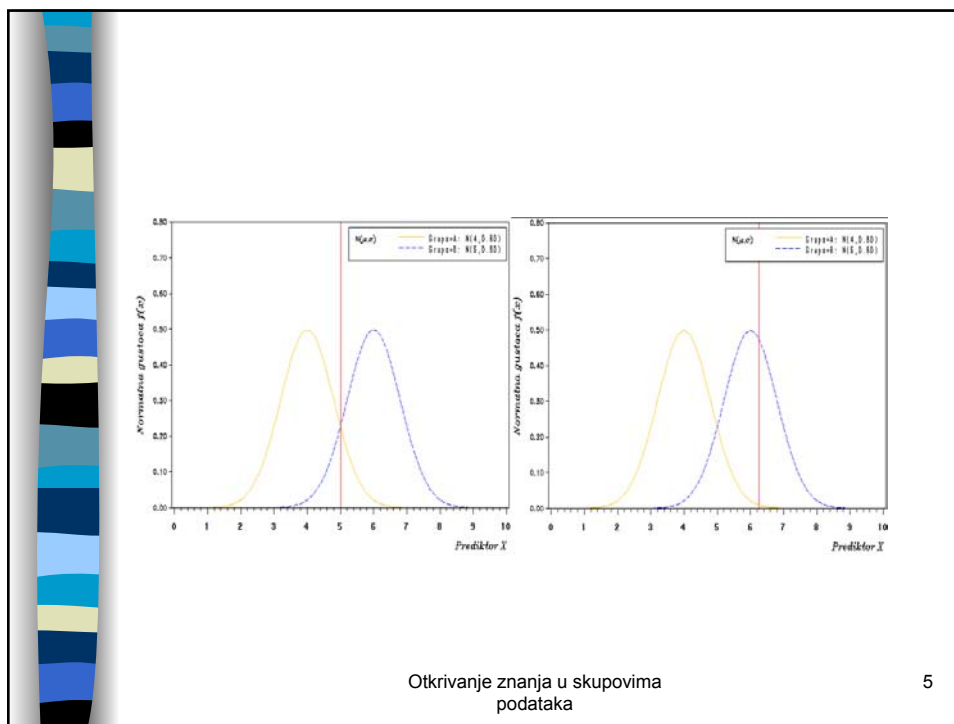
## Cilj diskriminativne analize

Naći funkcije prediktorskih varijabli koje  
maksimalno separiraju grupe opservacija

Klasifikacija novih opservacija u grupe

Otkrivanje znanja u skupovima  
podataka

4



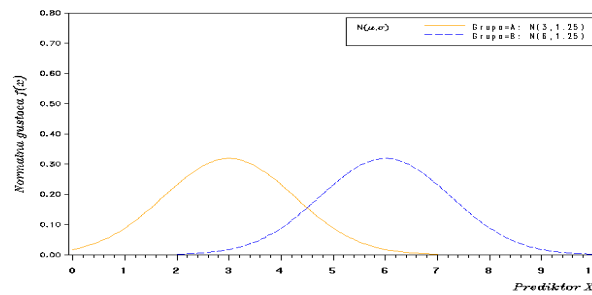
## Osnovni pojmovi

- opservacija: objekt iz populacije (grupe)
- **prediktorska varijabla** (nezavisna, prediktor):
  - kvantitativna, kontinuirana varijabla
  - njene su vrijednosti dobivene mjerenjem objekata
  - temeljem vrijednosti vrši se diskriminacija
- **kriterijska varijabla** (zavisna varijabla, kriterij, varijabla grupe):
  - kvalitativna, diskretna varijabla
  - određuje pripadnost objekta grupi
- separacija
- klasifikacija

Otkrivanje znanja u skupovima podataka

## Potrebni preduvjeti

1. Prediktorske varijable moraju biti multivarijatno normalno distribuirane u populacijama  $\Pi_i$  definiranim varijablom grupe  $i=1,2$

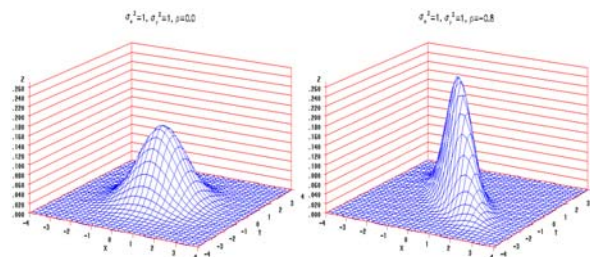


Otkrivanje znanja u skupovima podataka

7

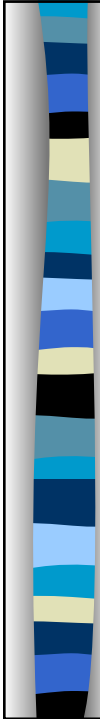
## Potrebni preduvjeti

2. Matrice varijanci-kovarijanci moraju biti homogene u populacijama  $\Pi_i$  definiranim varijablom grupe.



Otkrivanje znanja u skupovima podataka

8

- 
- Ako su uvjeti ispunjeni primjenjuje se **linearna diskriminativna funkcija**
  - Ako uvjeti nisu ispunjeni primjenjuje se kvadratna ili neparametarska diskriminativna funkcija



## Pravila za diskriminaciju grupa

- Pravilo vjerodostojnosti
  - funkcija gustoće
- Pravilo linearne diskriminativne funkcije
  - Fisherovo pravilo
- Pravilo Mahalanobisove udaljenosti
  - opisuje kvadratnu udaljenost opservacije od pripadne sredine grupe



## Broj diskriminativnih funkcija

- broj diskriminativnih funkcija za separaciju  $g$  grupa,  $g \geq 2$ , na osnovu  $p$  prediktorskih varijabli, jednak je:

$$\min(g-1, p)$$

*Iris data set*  $\min(3-1, 4) = 2$



## Temeljna ideja

Maksimizirati  
(SS između grupa)/(SS unutar grupa)



## Pravilo linearne diskriminativne funkcije

- Osnovni Fisherov cilj je separacija grupa (primjer dvije populacije) hiperravnina

$$Y = \mathbf{b}^T \mathbf{X} = \mathbf{b}_1 \mathbf{x}_1 + \mathbf{b}_2 \mathbf{x}_2 + \dots + \mathbf{b}_p \mathbf{x}_p = (\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)})^T \boldsymbol{\Sigma}^{-1} \mathbf{X}$$

Cilj je linearne funkcije  $Y$  maksimizirati (kvadriranu) udaljenost  $\boldsymbol{\mu}_Y^{(1)}$  i  $\boldsymbol{\mu}_Y^{(2)}$  u odnosu na varijabilitet od  $Y$ .

- klasifikacija

$$\mathbf{b}^T \mathbf{x} - k > 0$$

$$\text{uz } k = 1/2 (\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)})$$



## Testiranje pouzdanosti klasifikacije

Analiza pogrešaka klasifikacije po tipu i broju, uz vrednovanje različitog značenja pogrešaka

## Odabir značajnih prediktora

### Pitanje:

- da li su za efikasnu diskriminaciju potrebni svi prediktori
- koji su prediktori najbolji

Testira se statistička značajnost pojedine prediktorske varijable (statistički testovi);

### metode:

- selekcija unaprijed
- selekcija unatrag
- stepwise selekcija

## Primjer (Iris dataset)

### ■ Kriterijska varijabla (3 grupe)

- Iris Setosa, Iris Versicolor, Iris Virginica



### ■ Prediktorske varijable (4)

- SepalLength, SepalWidth,  
PetalLength, PetalWidth



Test (Kolmogorovljev D) normaliteta prediktora po grupama irisa

Prediktorske varijable	Vrsta irisa (grupa)		
	Setosa	Versicolor	Virginica
petal length	< 0.01	0.08	0.11
petal width	< 0.01	< 0.01	0.07
sepal length	0.10	> 0.15	0.09
sepal width	> 0.15	0.07	0.04

Otkrivanje znanja u skupovima podataka

17

Uz pretpostavku da su uvjeti diskriminativne analize zadovoljeni koeficijenti za formiranje linearnih funkcija bili bi:

$$\mathbf{b}_1 = [-0.083, -0.15, 0.22, 0.28]$$

$$\mathbf{b}_2 = [-0.002, 0.21, -0.09, 0.28]$$

Objе funkcije su statistički značajne pri čemu je prva statistički značajnija od druge (tj. prva funkcija jače diskriminira grupe od druge)

Otkrivanje znanja u skupovima podataka

18

## Koeficijenti za formiranje kanoničkih funkcija

Raw Canonical Coefficients		
Variable	Can1	Can2
SepalLength	-0.082	0.002
SepalWidth	-0.15	0.21
PetalLength	0.22	-0.09
PetalWidth	0.28	0.28

Otkrivanje znanja u skupovima podataka

19

## Provjera točnosti

Rezultati pokazuju da je krivo klasificirano 3% primjeraka (metoda krosvalidacije) odnosno 2% primjeraka (metoda resupstitucije)

## Odabir značajnih prediktora za diskriminaciju

- pokazalo se da su prediktori PetalLength, SepalWidth, PetalWidth, statistički značajni dok je SepalLength na granici značajnosti

Otkrivanje znanja u skupovima podataka

20

## Rezime rezultata metode u koracima (stepwise)

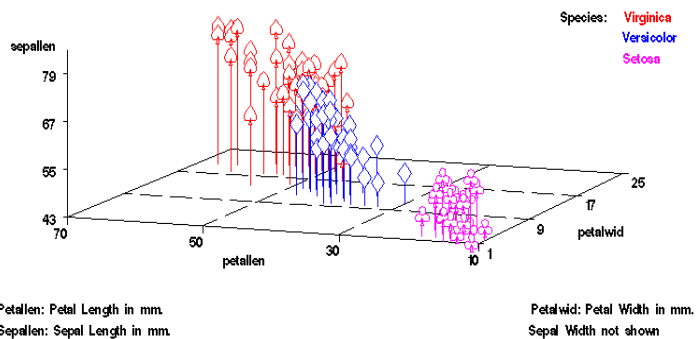
Stepwise Selection Summary								
Step	Number In	Entered	Removed	Partial R-Square	F Value	Pr > F	Wilks' Lambda	Pr < Lambda
1	1	PetalLength		0.94	1180.16	<.0001	0.05	<.0001
2	2	SepalWidth		0.37	43.04	<.0001	0.03	<.0001
3	3	PetalWidth		0.32	34.57	<.0001	0.02	<.0001
4	4	SepalLength		0.06	4.72	0.0103	0.02	<.0001

Otkrivanje znanja u skupovima podataka

21

## Iris Species Classification

Physical Measurement  
Source: Fisher (1936) Iris Data



Otkrivanje znanja u skupovima podataka

22



## Literatura

Sharma, S. (1996), "Applied Multivariate Techniques", John Wiley & Sons, Inc., pp. 237 – 316.

Johnson, R. A., Wichern, D. W. (1988), "Applied Multivariate Statistical Analysis", Prentice Hall, International Edition, pp. 470 – 542.