

Facultatea de Electronică, Telecomunicații și Tehnologie Informației

AIM AI Multimedia Lab
https://www.aimultimedialab.ro/

Universitatea Politehnică din București

Tehnici de analiză și clasificare automată a informației

Prof. dr. ing. Bogdan IONESCU
<https://bionescu.aimultimedialab.ro/>

București, 2022

1

Plan Curs

- M1. Introducere (concept, aplicații)
- M2. Prelucrarea și reprezentarea datelor de intrare
- M3. Tehnici de clasificare ne-supervizată ("clustering")
- M4. Tehnici de clasificare supervizată ("classification")
- M5. Evaluarea performanței clasificatorilor

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

2

> M3. Tehnici de clasificare ne-supervizată

- 3.1. [Introducere]
- 3.2. [Analiza similarității datelor]
- 3.3. [Clasificarea ierarhică]
- 3.4. [k-means]
- 3.5. [Gaussian Mixture Models]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

3

Clasificare ne-supervizată (clustering) - principiu

clustering = partiționarea datelor de intrare în mulțimi similare fără a dispune de informații a priori despre acestea (date de antrenare);

date de intrare

clasa 1

clasa 2

clasa 3

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

4

Clasificare ne-supervizată (clustering) - principiu (cont.)

clustering = partiționarea datelor de intrare în clase fără a dispune de exemple de partiționări (cont. exemplu);

date de intrare

clasa 1

clasa 2

clasa 3

clasa 4

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

5

Clasificare ne-supervizată (clustering) - principiu (cont.)

Întreg procesul depinde de modul de definire al conceptului de **similaritate** între date;

- similaritate = un concept foarte subiectiv;
- la nivel uman, este greu de definit dar îl recunoaștem atunci când îl vedem;

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

6

Analiza similarității datelor**1. Similaritatea descriptorilor**

determinarea gradului de asemănare dintre doi descriptori;

$$\left. \begin{array}{l} X_1 = [x_{1,1}, x_{1,2}, \dots, x_{1,n}] \\ \dots \\ X_m = [x_{m,1}, x_{m,2}, \dots, x_{m,n}] \end{array} \right\} d(X_i, X_j) = ?$$

instanțe de intrare

Dacă $d(\cdot)$ este metrică presupune:

- simetrie: $d(X_i, X_j) = d(X_j, X_i)$
- valoare minimă (0): $d(X_i, X_i)$
- respectă: $d(X_i, X_k) \leq d(X_i, X_j) + d(X_j, X_k) \quad \forall i, j, k$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

7

7

Analiza similarității datelor (cont.)**1. Similaritatea descriptorilor (cont.)**

↳ distanța Minkovski

$$d_{Mink}(X_i, X_j) = \sqrt[r]{\sum_{k=1}^n |x_{i,k} - x_{j,k}|^r}$$

unde X_i și X_j sunt două instanțe de intrare, $x_{i,k}$ cu $k=1, \dots, n$ reprezintă valorile atributelor pentru instanța X_i iar $|\cdot|$ reprezintă modulul.

↳ distanța Manhattan ($r=1$)

$$d_{Manh}(X_i, X_j) = \sum_{k=1}^n |x_{i,k} - x_{j,k}|$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

8

8

Analiza similarității datelor (cont.)**1. Similaritatea descriptorilor (cont.)**

↳ distanța Euclidiană ($r=2$)

$$d_{Euclid}(X_i, X_j) = \sqrt{\sum_{k=1}^n |x_{i,k} - x_{j,k}|^2}$$

> pentru a evita dependența de unitatea de măsură, datele pot fi standardizate = fiecare atribut să aibă pondere ~egală:

$$d_{wEuclid}(X_i, X_j) = \sqrt{\sum_{k=1}^n w_k \cdot |x_{i,k} - x_{j,k}|^2}$$

unde w_k sunt o serie de ponderi.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

9

9

Analiza similarității datelor (cont.)**1. Similaritatea descriptorilor (cont.)**

↳ distanța Canberra

$$d_{Canb}(X_i, X_j) = \frac{\sum_{k=1}^n |x_{i,k} - x_{j,k}|}{\sum_{k=1}^n (|x_{i,k}| + |x_{j,k}|)}$$

↳ distanța Bray-Curtis

$$d_{B-C}(X_i, X_j) = \frac{\sum_{k=1}^n |x_{i,k} - x_{j,k}|}{\sum_{k=1}^n (x_{i,k} + x_{j,k})}$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

10

10

Analiza similarității datelor (cont.)**1. Similaritatea descriptorilor (cont.)**

↳ distanța între date binare (valori 0 sau 1)

$$d_{bin}(X_i, X_j) = \frac{r + s}{q + r + s + t}$$

unde:

- q este numărul de atribute ce au valoarea 1 pentru ambele instanțe,
- t este numărul de atribute cu valoare 0 pentru ambele instanțe,
- $s + r$ reprezintă numărul de atribute de valori diferite pentru cele două instanțe (0 vs. 1 și respectiv 1 vs. 0).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

11

11

Analiza similarității datelor (cont.)**1. Similaritatea descriptorilor (cont.)**

↳ distanța între histogramme de valori

$$d_{inter}(X_i, X_j) = \sum_{k=1}^n \min\{x_{i,k}, x_{j,k}\}$$

unde $x_{i,k}$ cu $k=1, \dots, n$ (bini) reprezintă valorile histogrammei iar $\min\{\cdot\}$ returnează valoarea minimă a unei mulțimi.

$$d_{hist}(X_i, X_j) = \sqrt{(X_i - X_j)^T \cdot A \cdot (X_i - X_j)}$$

unde X reprezintă o histogramă, T este operația de transpusă iar $A = [a_{kl}]$ cu $k, l = 1, \dots, n$ este o matrice pătratică ce indică corelația dintre binii k și l .

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

12

12

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Bhattacharyya (între distribuții de probabilitate)

$$d_{\text{Bhatta}}(X_i, X_j) = \frac{1}{8}(\mu_{X_i} - \mu_{X_j})^T \cdot (\Sigma_{X_i, X_j})^{-1} \cdot (\mu_{X_i} - \mu_{X_j}) + \frac{1}{2} \cdot \ln \left(\frac{\det(\Sigma_{X_i, X_j})}{\sqrt{\det(\Sigma_{X_i}) \cdot \det(\Sigma_{X_j})}} \right)$$

unde μ_x este vectorul medie al distribuției de probabilitate a instanței X , Σ_x este matricea de covarianță a distribuției lui X , Σ_{X_i, X_j} este media aritmetică a matricelor de covarianță pentru distribuțiile lui X_i și X_j , T este transpusa unei matrice iar $\det(\cdot)$ returnează determinantul unei matrice.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 13

13

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Hausdorff (între mulțimi de valori)

$$d_{\text{Haus}}(X_i, X_j) = \max \left\{ \sup_k \left(\inf_l (d(x_{i,k}, x_{j,l})) \right), \sup_l \left(\inf_k (d(x_{i,k}, x_{j,l})) \right) \right\}$$

unde:

- $k, l = 1, \dots, n$;
- $\inf(\cdot)$ și $\sup(\cdot)$ sunt infimum și respectiv supremum al unei mulțimi;
- $d(\cdot)$ este o metrică;
- $\max\{\cdot\}$ returnează valoarea maximă a unei mulțimi.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 14

14

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Hausdorff (între mulțimi de valori; cont.)

$$d_{\text{Haus}}(X_i, X_j) = \max \left\{ \sup_k \left(\inf_l (d(x_{i,k}, x_{j,l})) \right), \sup_l \left(\inf_k (d(x_{i,k}, x_{j,l})) \right) \right\}$$

unde:

- $k, l = 1, \dots, n$;
- $\inf(\cdot)$ și $\sup(\cdot)$ sunt infimum și respectiv supremum al unei mulțimi;
- $d(\cdot)$ este o metrică;
- $\max\{\cdot\}$ returnează valoarea maximă a unei mulțimi.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 15

15

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Hausdorff (între mulțimi de valori; cont.)

$$d_{\text{Haus}}(X_i, X_j) = \max \left\{ \sup_k \left(\inf_l (d(x_{i,k}, x_{j,l})) \right), \sup_l \left(\inf_k (d(x_{i,k}, x_{j,l})) \right) \right\}$$

unde:

- $k, l = 1, \dots, n$;
- $\inf(\cdot)$ și $\sup(\cdot)$ sunt infimum și respectiv supremum al unei mulțimi;
- $d(\cdot)$ este o metrică;
- $\max\{\cdot\}$ returnează valoarea maximă a unei mulțimi.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 16

16

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Hausdorff (între mulțimi de valori; cont.)

$$d_{\text{Haus}}(X_i, X_j) = \max \left\{ \sup_k \left(\inf_l (d(x_{i,k}, x_{j,l})) \right), \sup_l \left(\inf_k (d(x_{i,k}, x_{j,l})) \right) \right\}$$

unde:

- $k, l = 1, \dots, n$;
- $\inf(\cdot)$ și $\sup(\cdot)$ sunt infimum și respectiv supremum al unei mulțimi;
- $d(\cdot)$ este o metrică;
- $\max\{\cdot\}$ returnează valoarea maximă a unei mulțimi.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 17

17

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Hausdorff (între mulțimi de valori; cont.)

$$d_{\text{Haus}}(X_i, X_j) = \max \left\{ \sup_k \left(\inf_l (d(x_{i,k}, x_{j,l})) \right), \sup_l \left(\inf_k (d(x_{i,k}, x_{j,l})) \right) \right\}$$

unde:

- $k, l = 1, \dots, n$;
- $\inf(\cdot)$ și $\sup(\cdot)$ sunt infimum și respectiv supremum al unei mulțimi;
- $d(\cdot)$ este o metrică;
- $\max\{\cdot\}$ returnează valoarea maximă a unei mulțimi.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 18

18

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța cosinus

$$d_{\cos}(X_i, X_j) = \frac{X_i \cdot X_j}{\|X_i\| \cdot \|X_j\|}$$

unde \cdot reprezintă produsul scalar iar $\|\cdot\|$ reprezintă norma unui vector, astfel:

$$\|X\|^2 = \sum_{k=1}^n x_k^2$$

unde $X = [x_1, x_2, \dots, x_n]$.

> distanța este practic cosinusul unghiului celor doi vectori normalizați.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 19

19

Analiza similarității datelor (cont.)

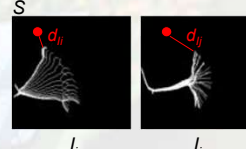
1. Similaritatea descriptorilor (cont.)

↳ distanța Baddeley (între obiecte)

$$d_{\text{Badd}}(I_i, I_j) = \left[\frac{1}{M \cdot N} \sum_{p \in S} |d_{I_i}(p) - d_{I_j}(p)|^q \right]^{\frac{1}{q}}$$

unde:

- I este o imagine binară,
- S reprezintă setul de puncte din imagine ($M \times N$ puncte),
- $d_i(p)$ reprezintă o anumită metrică de la punctul p la cel mai apropiat punct al obiectului din imaginea I ,
- q este exponentul (ex. $q=2$).



Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 20

20

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Earth Mover's Distance (între date de dimensiuni diferite)

$$d_{\text{EMD}}(X_i, X_j) = \frac{\sum_{k=1}^m \sum_{l=1}^n d_{k,l} \cdot f_{k,l}}{\sum_{k=1}^m \sum_{l=1}^n f_{k,l}}$$

unde X_i și X_j au dimensiuni diferite (m, n), $d_{k,l}$ reprezintă distanța dintre valorile $x_{i,k}$ și $x_{j,l}$ iar $f_{k,l}$ este o funcție de cost ce cuantizează deplasarea între $x_{i,k}$ și $x_{j,l}$ determinată ca minimizând valoarea costului total:

$$\sum_{k=1}^m \sum_{l=1}^n d_{k,l} \cdot f_{k,l}$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 21

21

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Earth Mover's Distance (cont.)

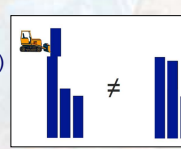
> reprezintă practic "volumul de muncă" necesar transformării unei instanțe în cealaltă;

> exemplu, fie:

$$X = [(x_1, w_1), (x_2, w_2), \dots, (x_n, w_n)]$$

$$Y = [(y_1, u_1), (y_2, u_2), \dots, (y_n, u_n)]$$

unde X și Y sunt două instanțe de comparat iar w și u reprezintă ponderile atributelor (~masă);



Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 22

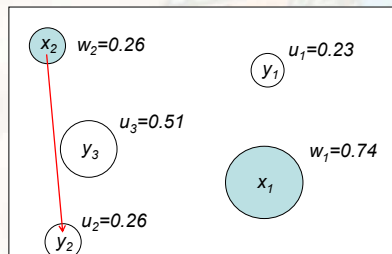
22

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Earth Mover's Distance (cont.)

> exemplu (cont.) - presupunem ponderi egale, $\Sigma w = \Sigma u$;



- calculăm necesarul de muncă ca să transformăm X în Y (mutăm masa de la X la Y);

$$\text{work}_{2,2} = f_{2,2} \cdot d_{2,2} = 0.26 \cdot 316$$

[S. Cohen, 1999]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 23

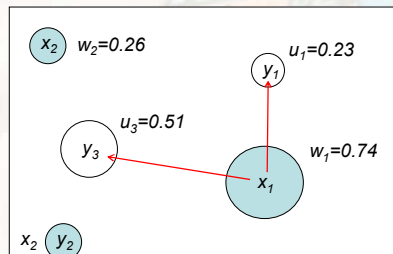
23

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Earth Mover's Distance (cont.)

> exemplu (cont.) - presupunem ponderi egale, $\Sigma w = \Sigma u$;



$$\text{work}_{1,1} = f_{1,1} \cdot d_{1,1} = 0.23 \cdot 155$$

$$\text{work}_{1,3} = f_{1,3} \cdot d_{1,3} = 0.51 \cdot 252$$

[S. Cohen, 1999]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 24

24

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

↳ distanța Earth Mover's Distance (cont.)

> exemplu (cont.) - presupunem ponderi egale, $\Sigma w = \Sigma u$;

$$twork = 0.23 * 155 + 0.51 * 252 + 0.26 * 316 = 246$$

Este corect ca măsură de distanță ?

Nu au fost alese costurile optime!

[S. Cohen, 1999]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 25

25

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

$twork = 222$

↳ distanța Earth Mover's Distance (cont.)

> exemplu (cont.) - presupunem ponderi egale, $\Sigma w = \Sigma u$;

- optimal:

$work_{2,3} = f_{2,3} * d_{2,1} = 0.26 * 198$

$work_{1,1} = f_{1,1} * d_{1,1} = 0.23 * 155$

$work_{1,2} = f_{1,2} * d_{1,2} = 0.26 * 277$

$work_{1,3} = f_{1,3} * d_{1,3} = 0.25 * 252$

[S. Cohen, 1999]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 26

26

Analiza similarității datelor (cont.)

1. Similaritatea descriptorilor (cont.)

> cât de mult contează alegerea adecvată a măsurii de distanță adaptată datelor?

măsură de performanță (valoare > performanță > maxim 1 - 100%)

[I. Mironică et al., EUSIPCO 2012]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 27

27

Analiza similarității datelor (cont.)

2. Similaritatea la nivel de structură

determinarea gradului de asemănare a două obiecte la nivel structural (ex. aranjare spațială, structurare text, etc);

> exemplu, compararea a două documente video;

> idee, reprezentare textuală a structurii temporale (vezi M2):

- "s" – plan video;
- "c" – tranziție de tip cut;
- "w" – tranziție de tip wipe;
- "d" – tranziție de tip dissolves;

> documentul video este reprezentat ca o secvență de litere:

"scswsdcs"

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 28

28

Analiza similarității datelor (cont.)

2. Similaritatea la nivel de structură (cont.)

> exemplu, compararea a două documente video (cont.);

↳ distanța de editare

costul minim de transformare a instanței X_i în instanța X_j , unde X_i și X_j au n și respectiv m caractere ce pot lua valori într-un alfabet Σ iar E definește setul de operații de editare și costurile acestora.

$$\left. \begin{matrix} X_i = \text{"scswsdcs"} \\ X_j = \text{"sdswscscs"} \end{matrix} \right\} d(X_i, X_j) = \begin{matrix} 2 \text{ înlocuiri} + \\ 1 \text{ inserare} \end{matrix} = 1 + 1 + 1 = 3$$

$\Sigma = \{c, w, d, s\}$

$E = \{\text{"inserare"}, \text{"ștergere"}, \text{"înlocuire"}\}$ (costuri egale, 1)

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 29

29

Analiza similarității datelor (cont.)

3. Similaritatea semantică

determinarea gradului de asemănare la nivel de concepte (reprezentare semantică a informației);

> ontologii de informații:

- mod formal de reprezentare a cunoașterii sub formă de concepte și a relațiilor dintre acestea;
- folosesc următoarele componente:
 - obiecte/instanțe de date;
 - clase (mulțimi, colecții, concepte);
 - atribute (proprietăți);
 - relații (între clase și instanțe);
 - restricții;
 - reguli (de tip "if-then");
 - evenimente (modul de schimbare al atributelor).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 30

30

Analiza similarității datelor (cont.)

3. Similaritatea semantică (cont.)

> ontologii de informații (cont.):

ontologie clasă "mașină" (simplificat)

- structură ierarhică de reprezentare a informației;
- clase subordonate moștenesc proprietățile claselor superioare;
- obiectele sunt descrise de atribute, exemplu:
 - nume "Ford Explorer";
 - ușă (4);
 - motor (4l);
 - transmisie (6v).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 31

31

Analiza similarității datelor (cont.)

3. Similaritatea semantică (cont.)

> ontologii de informații (cont.):

↳ distanța între concepte

ontologie

- datele sunt reprezentate pe bază de ontologii semantice;
- exemplu: distanța dintre instanța c_1 și respectiv c_2 (descriptori) = numărul de pași din arbore necesari pentru a ajunge de la conceptul C_1 (definește c_1) la conceptul C_2 (definește c_2);
- număr de pași = număr de laturi (3 în exemplu).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 32

32

Clasificarea ierarhică (hierarchical clustering)

datele de intrare sunt regrupate în funcție de similaritatea acestora într-un număr variabil de clase (1-n), sub formă arborescentă;

[sursă imagine Wikipedia]

- complexitate de calcul redusă;
- numărul de clase rezultate poate fi adaptat în funcție de scenariu;
- aglomerativ sau "**bottom up**" – de jos în sus (în figura de alături);
- diviziv sau "**top down**" – de sus în jos (în figura de alături).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 33

33

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă

> date de intrare, $X_i = [x_{i,1}, \dots, x_{i,m}]$, $i = 1, \dots, m$;

> algoritm:

p1. fiecare dintre instanțe este asociată unei clase
 -> $X_i \in \text{clasa}_1, \dots, X_m \in \text{clasa}_m$;

p2. se calculează pentru fiecare pereche de clase o măsură de similaritate între acestea;

	clasa ₁	clasa ₂	...	clasa _m
clasa ₁	0	$d(1,2)$...	$d(1,m)$
...	$d(2,1)$	0	...	$d(2,m)$
clasa _m	0

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 34

34

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> algoritm (cont.):

p3. cele mai similare două clase sunt fuzionate într-una singură;

	clasa _{1&2}	clasa ₃	...	clasa _m
clasa _{1&2}	0	$d(1\&2,3)$...	$d(1\&2,m)$
...	$d(3,1\&2)$	0	...	$d(3,m)$
clasa _m	0

p4. dacă numărul de clase obținute ≤ 1 , mergi la pasul 2 (se re-calculează similaritatea între noile clase și se fuzionează în continuare);

p5. STOP -> dendrograma claselor.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 35

35

Clasificarea ierarhică (hierarchical clustering; cont.)

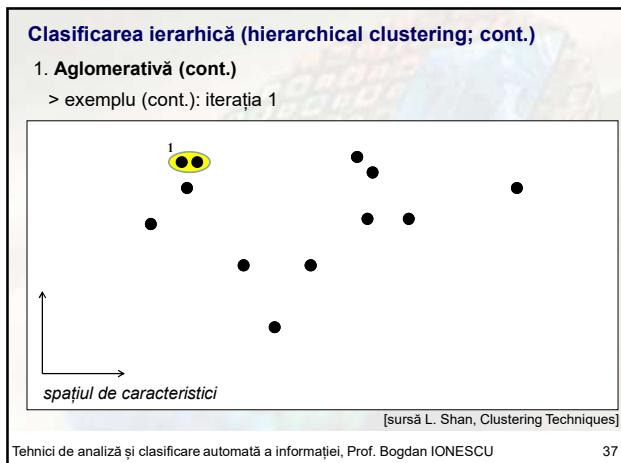
1. Aglomerativă (cont.)

> exemplu:

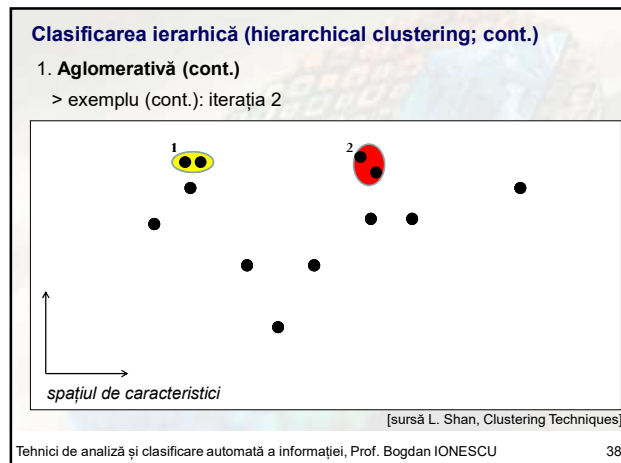
[sursă L. Shan, Clustering Techniques]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 36

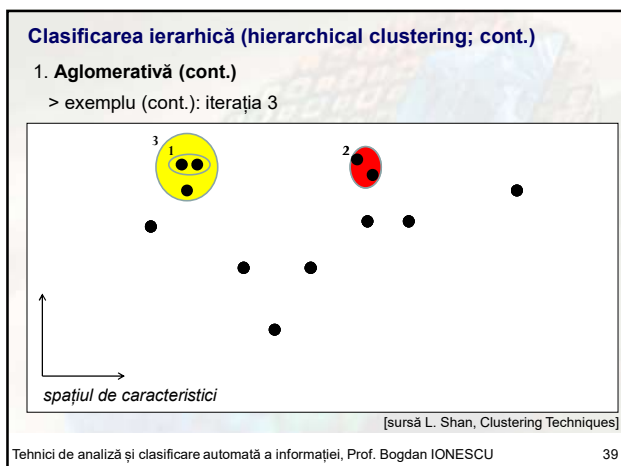
36



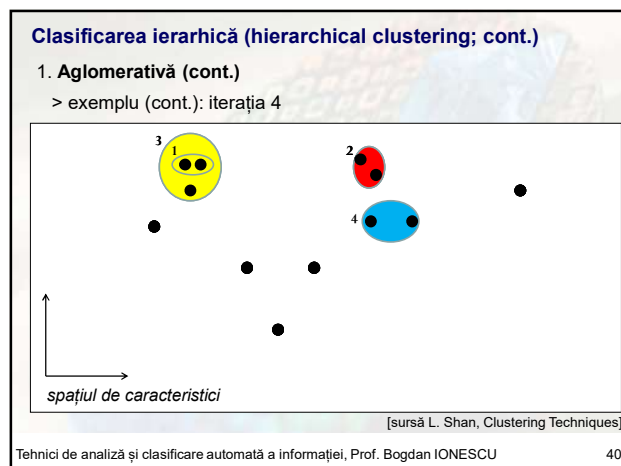
37



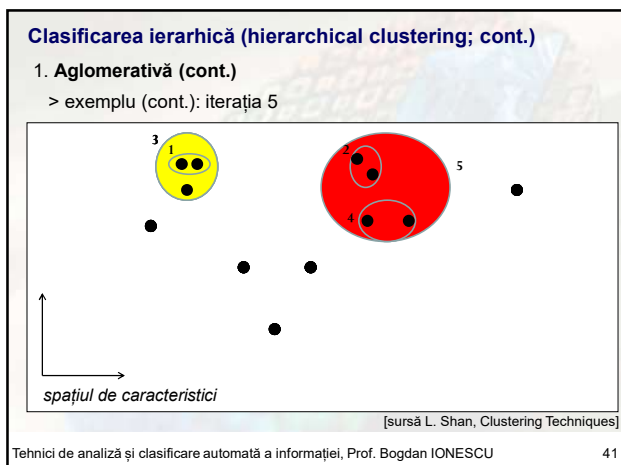
38



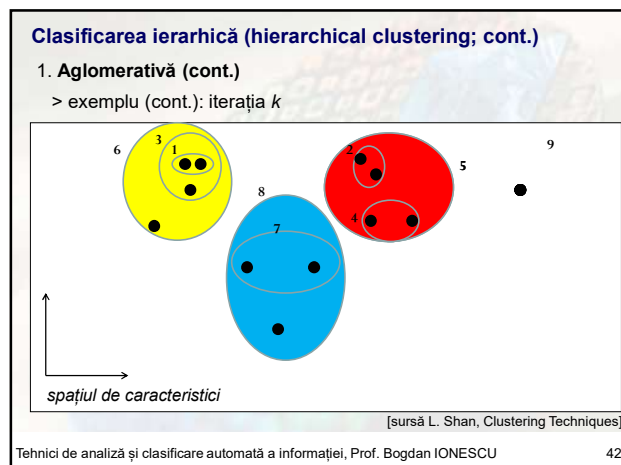
39



40



41



42

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> exemplu (cont.): iterația 11

spațiul de caracteristici

[sursă L. Shan, Clustering Techniques]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 43

43

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> cum evaluăm similaritatea între clase?

spațiul de caracteristici

> *single link* = distanța dintre cele mai apropiate două instanțe ale claselor;

-> clasele rezultate tind să fie subțiri și lungi.

[sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 44

44

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> cum evaluăm similaritatea între clase? (cont.)

spațiul de caracteristici

> *complete link* = distanța dintre cele mai depărtate două instanțe ale claselor;

-> clasele rezultate tind să fie foarte apropiate.

[sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 45

45

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> cum evaluăm similaritatea între clase? (cont.)

spațiul de caracteristici

> *average link* = distanța medie dintre toate instanțele celor două clase;

-> robustețe la zgomot.

[sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 46

46

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link): 10 clase

[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Battlo"]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 47

47

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 9 clase

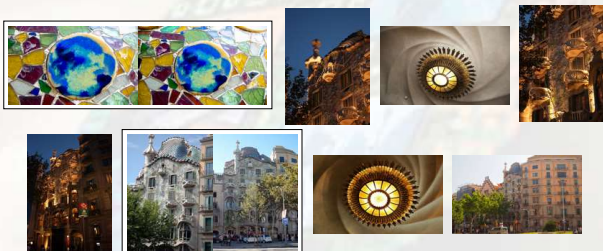
[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Battlo"]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 48

48

Clasificarea ierarhică (hierarchical clustering; cont.)**1. Aglomerativă (cont.)**

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 8 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batllo"]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

49

49

Clasificarea ierarhică (hierarchical clustering; cont.)**1. Aglomerativă (cont.)**

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 7 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batllo"]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

50

50

Clasificarea ierarhică (hierarchical clustering; cont.)**1. Aglomerativă (cont.)**

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 6 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batllo"]

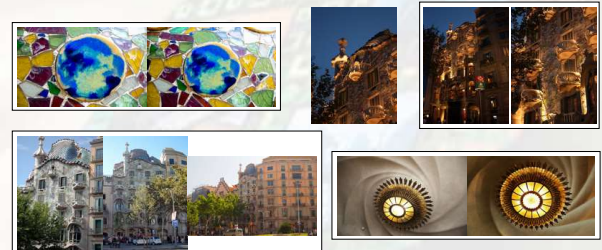
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

51

51

Clasificarea ierarhică (hierarchical clustering; cont.)**1. Aglomerativă (cont.)**

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 5 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batllo"]

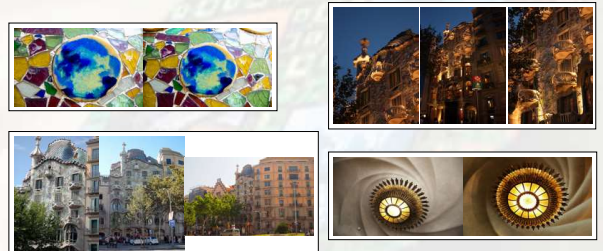
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

52

52

Clasificarea ierarhică (hierarchical clustering; cont.)**1. Aglomerativă (cont.)**

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 4 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batllo"]

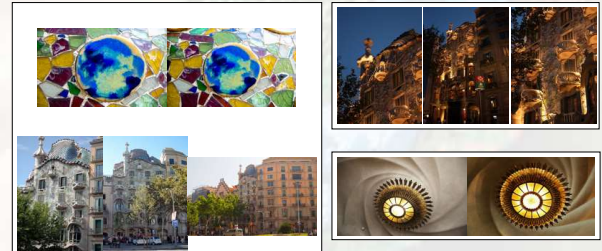
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

53

53

Clasificarea ierarhică (hierarchical clustering; cont.)**1. Aglomerativă (cont.)**

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 3 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batllo"]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

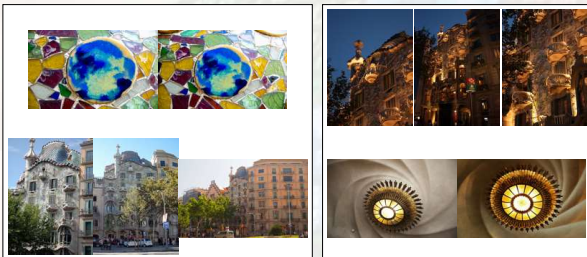
54

54

Clasificarea ierarhică (hierarchical clustering; cont.)

1. Aglomerativă (cont.)

> exemplu real (imagini, descriptor - culoare, metrică - Euclidiană, similaritate - average link; cont.) - 2 clase



[credit B. Boteanu, 2015; sursă imagini Flickr, "Casa Batlló"]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

55

55

Clasificarea ierarhică (hierarchical clustering; cont.)

2. Divizivă

> date de intrare, $X_i = [x_{i,1}, \dots, x_{i,m}]$, $i=1, \dots, m$;

> algoritm:

- p1. toate instanțele sunt asociate unei singure clase
-> $X_1, \dots, X_m \in \text{clasa}_1$;
- p2. clasele curente sunt divizate în două subclase folosind orice algoritm de partiționare;
- p3. dacă numărul de clase $< m$ se repetă pasul 2;
- p4. STOP -> dendrograma claselor.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

56

56

Clasificarea ierarhică (hierarchical clustering; cont.)

> care dintre cele două abordări ("top down" vs. "bottom up") este mai complexă?

= top down pentru că necesită o altă metodă de clustering;

> care dintre cele două abordări tinde să fie mai eficientă?

= top down, complexitate liniară funcție de numărul de clase (folosind k-means pentru partiționare);
vs. bottom up, cel puțin pătratică.

> care dintre cele două abordări tinde să fie mai precisă?

- bottom up – deciziile de agreare sunt luate local fără a ține cont de distribuția globală (deciziile inițiale nu mai pot fi schimbate ulterior);
- top down – țin cont de distribuția globală.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

57

57

k-means

partiționarea iterativă a datelor în k clase în funcție de proximitatea acestora față de reprezentanții claselor (centrozii);

> date de intrare:

- instanțe de clasificat în k clase:

$$X = \{X_1, X_2, \dots, X_m\} \rightarrow c_1, \dots, c_k;$$

- un dicționar de k instanțe:

$$V = \{V_1, V_2, \dots, V_k\}$$

- o matrice de partiționare:

$$\Gamma = [\gamma_{l,i}], \gamma_{l,i} = \begin{cases} 1 & X_i \in c_l \\ 0 & \text{altfel} \end{cases}$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

58

58

k-means (cont.)

> algoritm:

$$\text{optimizare } E(\Gamma, V) = \sum_{l=1}^k \sum_{i=1}^m \gamma_{l,i} \|X_i - V_l\|^2$$

- p1. se alege o valoare pentru k (numărul de clase);
- p2. se inițializează vocabularul V cu k instanțe din datele de intrare X . Acestea definesc o partiție inițială a claselor (centrozii);
- p3. fiecare instanță este atribuită clasei celei mai apropiate (ca distanță față de centroidul clasei);
- p4. se calculează matricea Γ de partiționare în clase;
- p5. se re-calculează vocabularul, fiecare vector fiind înlocuit de centroidul (= media) clasei respective;
- p6. se reia pasul 3 până când nici o instanță nu-și mai schimbă apartenența la clase (Γ nu se modifică).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

59

59

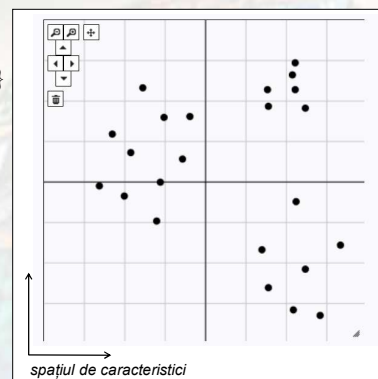
k-means (cont.)

> exemplu:

$$X = \{X_1, \dots, X_{23}\}$$

$$k = 3$$

$$c_1, c_2, c_3$$



Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU

60

60

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- se alege vocabularul din instanțele de intrare;
- acesta definește cele 3 clase;
- instanțele sunt asociate claselor cele mai apropiate.

spațiul de caracteristici

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 61

61

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- se recalculază vectorii V pentru fiecare clasă ca fiind centrozii claselor curente (medie);

spațiul de caracteristici

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 62

62

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- instanțele sunt asociate claselor celor mai apropiate pe baza noilor centrozii;

spațiul de caracteristici

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 63

63

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- se repetă pașii anteriori ...;

spațiul de caracteristici

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 64

64

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- etc;

spațiul de caracteristici

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 65

65

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- etc;

spațiul de caracteristici

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 66

66

k-means (cont.) [sursă <http://util.io/k-means>]

> exemplu (cont.):

- în acest moment nu se mai schimbă repartitia în clase a instanțelor;

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 67

67

k-means (cont.)

> avantaje:

- simplu de implementat;
- optimizează în mod intuitiv similaritatea intra-clasă;
- relativ eficient, complexitate $O(m \times k \times nr. \text{iterații})$.

> dezavantaje:

- necesită definirea noțiunii de centroid ca medie instanțe;
- optimizare locală – depinde practic de alegerea (bună) a vocabularului inițial pentru clase;
- numărul de clase trebuie anticipat;
- sensibil la date atipice;
- nu este eficient pentru clustere cu forme non-convexe.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 68

68

k-means (cont.)

> date atipice (outliers):

- potențială soluție: *k-medoids*, centrele claselor sunt alese ca fiind chiar unele dintre instanțe și nu mediile;

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 69

69

k-means (cont.)

> clustere cu forme non-convexe?

- potențială soluție: *kernel trick*;
- funcția de cost standard de minimizat este:

$$E(\Gamma, V) = \sum_{l=1}^k \sum_{i=1}^m \gamma_{li} \|X_i - V_l\|^2$$

- idee: transformăm X printr-o funcție (nucleu – kernel):

$$E(\Gamma, V) = \sum_{l=1}^k \sum_{i=1}^m \gamma_{li} \|\varphi(X_i) - \varphi(V_l)\|^2$$

$$\varphi(V_l) = \frac{1}{m_l} \sum_{i=1}^m \gamma_{li} \varphi(X_i), m_l \text{ este numărul de instanțe din clasa } l.$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 70

70

k-means (cont.)

> clustere cu forme non-convexe (cont.)

- potențială soluție: *kernel trick* (cont.);

$$E(\Gamma, V) = \sum_{l=1}^k \sum_{i=1}^m \gamma_{li} \|\varphi(X_i) - \varphi(V_l)\|^2$$

$$\|\varphi(X_i) - \varphi(V_l)\|^2 = \varphi(X_i)^T \cdot \varphi(X_i) - \varphi(X_i)^T \cdot \varphi(V_l) - \varphi(V_l)^T \cdot \varphi(X_i) + \varphi(V_l)^T \cdot \varphi(V_l)$$

- funcția nucleu este dată de: $\varphi(X_i)^T \cdot \varphi(X_j) = K(X_i, X_j)$
- $\|\varphi(X_i) - \varphi(V_l)\|^2 = K(X_i, X_i) - 2K(X_i, V_l) + K(V_l, V_l)$
- exemple de nuclee: $K(X_i, X_j) = e^{-\sigma \|X_i - X_j\|}$ (Gaussian);
- $K(X_i, X_j) = (c + X_i^T \cdot X_j)^d$ (polinomial).

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 71

71

k-means (cont.)

> clustere cu forme non-convexe (cont.)

- potențială soluție: *kernel trick* (cont.);

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 72

72

k-means (cont.)

- > cum putem decide în mod "automat" numărul de clase?
- > idee: pentru setul de date, încercăm mai multe valori pentru k :

spațiul de caracteristici [sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 73

73

k-means (cont.)

- > cum putem decide în mod "automat" numărul de clase (cont.)
- > idee: pentru setul de date, încercăm mai multe valori pentru k :

spațiul de caracteristici [sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 74

74

k-means (cont.)

- > cum putem decide în mod "automat" numărul de clase (cont.)
- > idee: pentru setul de date, încercăm mai multe valori pentru k :

spațiul de caracteristici [sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 75

75

k-means (cont.)

- > cum putem decide în mod "automat" numărul de clase (cont.)
- > idee: pentru setul de date, încercăm mai multe valori pentru k :

[sursă H. Lin, 15-381 Artificial Intelligence]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 76

76

Gaussian Mixture Models

abordare bazată pe modele; clasele sunt considerate a avea distribuții Gaussiene ale căror parametri sunt optimizați astfel încât să se potrivească cel mai bine datelor;

- funcția de repartiție:

$$F_X(x) = P\{X \leq x\}$$

unde X este o variabilă aleatoare, x reprezintă o valoare iar $P\{\}$ reprezintă probabilitatea în sensul statistic.

- > reprezintă probabilitatea ca realizarea particulară a variabilei aleatoare X să fie mai mică sau egală decât x .

$$0 \leq F_X(x) \leq 1$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 77

77

Gaussian Mixture Models (cont.)

- funcția de densitate de probabilitate:

$$f_X(x) = \frac{dF_X(x)}{dx}, f_X(x) \geq 0$$

unde d/dx reprezintă derivata de ordinul 1.

$$P\{X \leq x\} = F_X(x) = \int_{-\infty}^x f_X(t) dt$$

- > aria de sub graficul format de densitatea de probabilitate.

$$P\{x_1 < X \leq x_2\} = \int_{x_1}^{x_2} f_X(t) dt$$

$$P\{X \approx x\} = f_X(x) dx = P\{x < X \leq x + dx\}$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 78

78

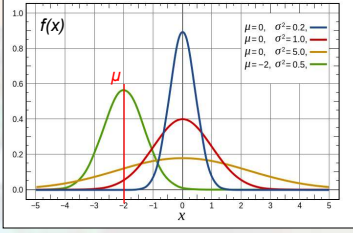
Gaussian Mixture Models (cont.) [sursă imagine Wikipedia]

- densitate de probabilitate normală, Gaussiană (1D):

$$N(X; \mu, \sigma^2) : f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

unde μ reprezintă media valorilor și σ este abaterea pătratică medie.

> 68% din valori sunt în intervalul $[\mu-\sigma; \mu+\sigma]$;
 > 99% din valori sunt în intervalul $[\mu-3\sigma; \mu+3\sigma]$.



Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 79

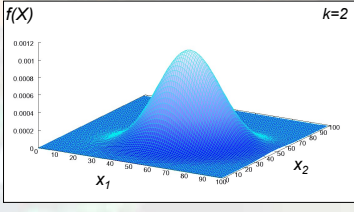
79

Gaussian Mixture Models (cont.) [sursă imagine Wikipedia]

- densitate de probabilitate normală, Gaussiană (nD):

$$N(X; \mu, \Sigma) : f(x_1, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} e^{-\frac{1}{2}(X-\mu)^T \Sigma^{-1} (X-\mu)}$$

unde $X=[x_1, \dots, x_k]$ reprezintă o variabilă aleatoare k dimensională,
 $\mu=[\mu_1, \dots, \mu_k]$ reprezintă vectorul medie (μ_i este media lui x_i), Σ este matricea de covarianță (dimensiune $k \times k$),
 T reprezintă transpusa,
 $^{-1}$ reprezintă inversa iar $\det()$ returnează determinantul.



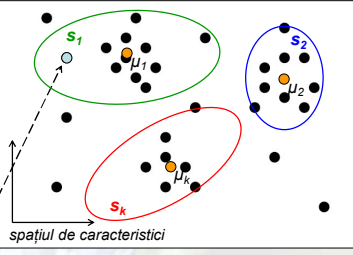
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 80

80

Gaussian Mixture Models (cont.)

> ipoteza GMM:

- se presupune faptul că avem la îndemână k surse de date;
- fiecare sursă i generează date de medie μ_i și matrice de covarianță Σ_i (distribuție Gaussiană);
- astfel, pentru o sursă i , de probabilitate p_i , datele generate de aceasta au distribuție $\sim N(\mu_i, \Sigma_i)$.



Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 81

81

Gaussian Mixture Models (cont.)

> clasificator: determinarea optimă a acestor distribuții ce se potrivesc cel mai bine repartiției datelor de intrare în spațiul de caracteristici (amestec de Gaussiene - GMM);

> optimizare = algoritm Expectation-Maximization (EM);

> date de intrare:

- instanțele de clasificat în k clase:
 $X = \{X_1, X_2, \dots, X_m\} \rightarrow c_1, \dots, c_k$;
- probabilitățile celor k surse:
 p_1, \dots, p_k
- valorile medii și matricele de covarianță:
 $\mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 82

82

Gaussian Mixture Models (cont.)

> algoritm GMM + EM:

p1. se alege numărul de surse k (= numărul de clase);

p2. se inițializează parametri de intrare, p_i, μ_i, Σ_i cu $i=1, \dots, k$ (ex. valori aleatorii);

$$\lambda = \{ \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k, p_1, \dots, p_k \}$$

p3. sunt calculate clasele estimate (Expectation-step):

$$P\{c_i | X_j, \lambda\} = \frac{P\{X_j | c_i, \lambda\} \cdot P\{c_i | \lambda\}}{P\{X_j | \lambda\}}$$

se eval. $N(X_j; \mu_i, \Sigma_i)$

$$P\{X_j | c_i, \lambda\} \cdot P\{c_i | \lambda\} = P\{X_j | c_i, \mu_i, \Sigma_i\} \cdot p_i$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 83

83

Gaussian Mixture Models (cont.)

> algoritm GMM + EM (cont.):

p3. sunt calculate clasele estimate (Expectation-step; cont):

$$P\{c_i | X_j, \lambda\} = \frac{P\{X_j | c_i, \lambda\} \cdot P\{c_i | \lambda\}}{P\{X_j | \lambda\}}$$

$$P\{X_j | c_i, \lambda\} \cdot P\{c_i | \lambda\} = P\{X_j | c_i, \mu_i, \Sigma_i\} \cdot p_i$$

$$P\{X_j | \lambda\} = \sum_{i=1}^k P\{X_j | c_i, \mu_i, \Sigma_i\} \cdot p_i$$

se eval. $N(X_j; \mu_i, \Sigma_i)$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 84

84

Gaussian Mixture Models (cont.)

> algoritm GMM + EM (cont.):

p4. sunt maximizate mediile și sunt recalculați parametrii (Maximization-step):

$$\mu_i = \frac{\sum_{j=1}^m P\{c_i | X_j, \lambda\} \cdot X_j}{\sum_{j=1}^m P\{c_i | X_j, \lambda\}}$$

$$\Sigma_i = \frac{\sum_{j=1}^m P\{c_i | X_j, \lambda\} [X_j - \mu_i][X_j - \mu_i]^T}{\sum_{j=1}^m P\{c_i | X_j, \lambda\}}$$

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 85

85

Gaussian Mixture Models (cont.)

> algoritm GMM + EM (cont.):

p4. sunt maximizate mediile și sunt recalculați parametrii (Maximization-step; cont.):

$$p_i = \frac{\sum_{j=1}^m P\{c_i | X_j, \lambda\}}{m}$$

p5. dacă parametrii de intrare, în urma actualizării, se schimbă foarte puțin -> STOP;

p6. altfel se repetă procesul cu pasul 3.

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 86

86

Gaussian Mixture Models (cont.)

> exemplu: iterația 0

- inițializare:
- k=3;
- probabilități surse egale (0.33);
- medii;
- matrice de covarianță.
- calcul probabilități de apartenență la distribuții;

spațiul de caracteristici [sursă Andrew W. Moore]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 87

87

Gaussian Mixture Models (cont.)

> exemplu: iterația 1

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...

spațiul de caracteristici [sursă Andrew W. Moore]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 88

88

Gaussian Mixture Models (cont.)

> exemplu: iterația 2

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...

spațiul de caracteristici [sursă Andrew W. Moore]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 89

89

Gaussian Mixture Models (cont.)

> exemplu: iterația 3

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...

spațiul de caracteristici [sursă Andrew W. Moore]

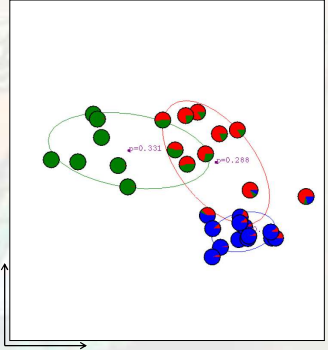
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 90

90

Gaussian Mixture Models (cont.)

> exemplu: iterația 4

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...



[sursă Andrew W. Moore]

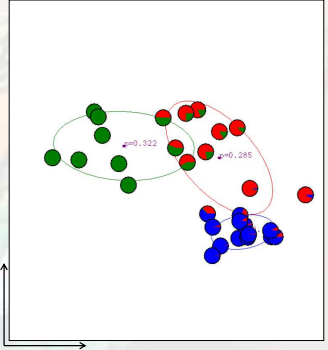
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 91

91

Gaussian Mixture Models (cont.)

> exemplu: iterația 5

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...



[sursă Andrew W. Moore]

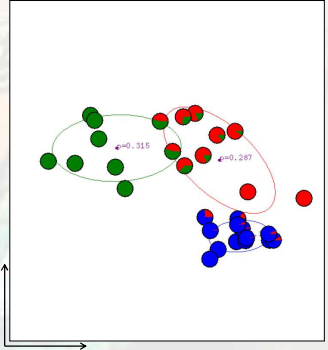
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 92

92

Gaussian Mixture Models (cont.)

> exemplu: iterația 6

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...



[sursă Andrew W. Moore]

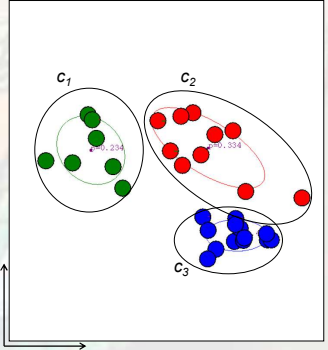
Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 93

93

Gaussian Mixture Models (cont.)

> exemplu: iterația 20

- în urma actualizării se redistribuie probabilitățile claselor, medii și matricea de covarianță ...
- rezultă repartiția optimă în clase de distribuție normală.



[sursă Andrew W. Moore]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 94

94

Gaussian Mixture Models (cont.)

> avantaje:

- interpretabilitate: determină un model de generare a datelor (se pot genera date noi);
- relativ eficient, complexitate $O(m \times k \times nr. \text{iterații})$;
- extensibil la alt tip de distribuții de date.

> dezavantaje:

- EM conduce de regulă la un minim local – depinde de inițializare;
- numărul de clase trebuie determinat a priori;
- mai puțin eficient pentru clase de formă ne convexă.

[sursă Andrew W. Moore]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 95

95

> Sfârșit M3

[sursă Andrew W. Moore]

Tehnici de analiză și clasificare automată a informației, Prof. Bogdan IONESCU 96

96