

Învățarea nesupervizată Clustering

Învățarea automată

Direcții principale

1. Clasificarea și regresia
2. Gruparea (clustering [clástering], clusterizare [clasterizáre])
3. Determinarea regulilor de asociere
4. Selecția trăsăturilor

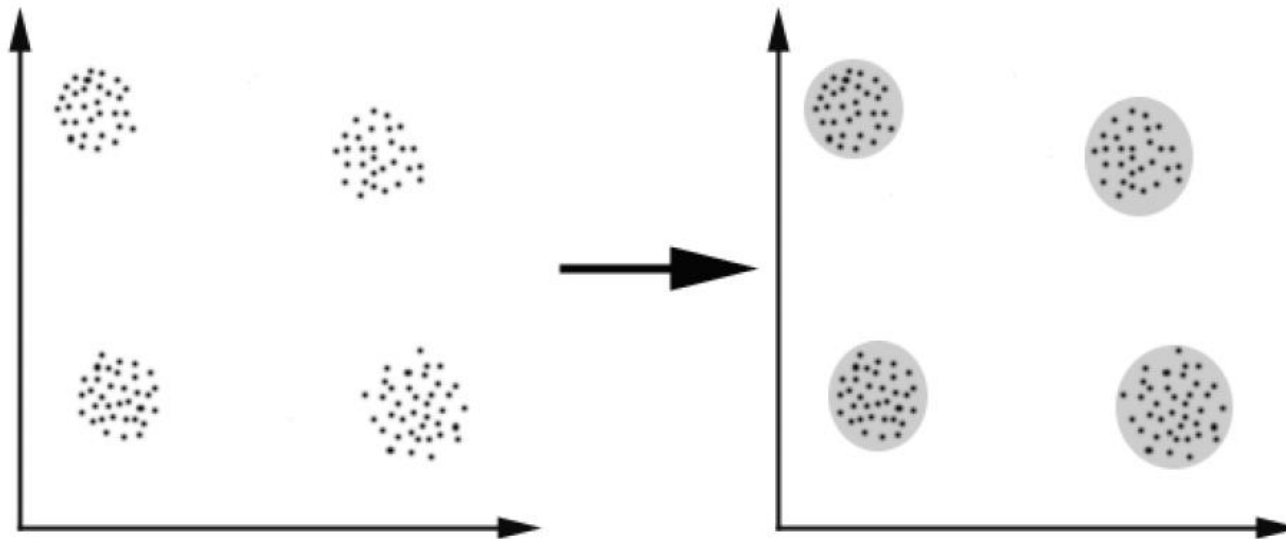
Învățarea automată

Direcții principale

1. Clasificarea și regresia
2. Gruparea (clustering [clástering], clusterizare [clasterizáre])
3. Determinarea regulilor de asociere
4. Selecția trăsăturilor

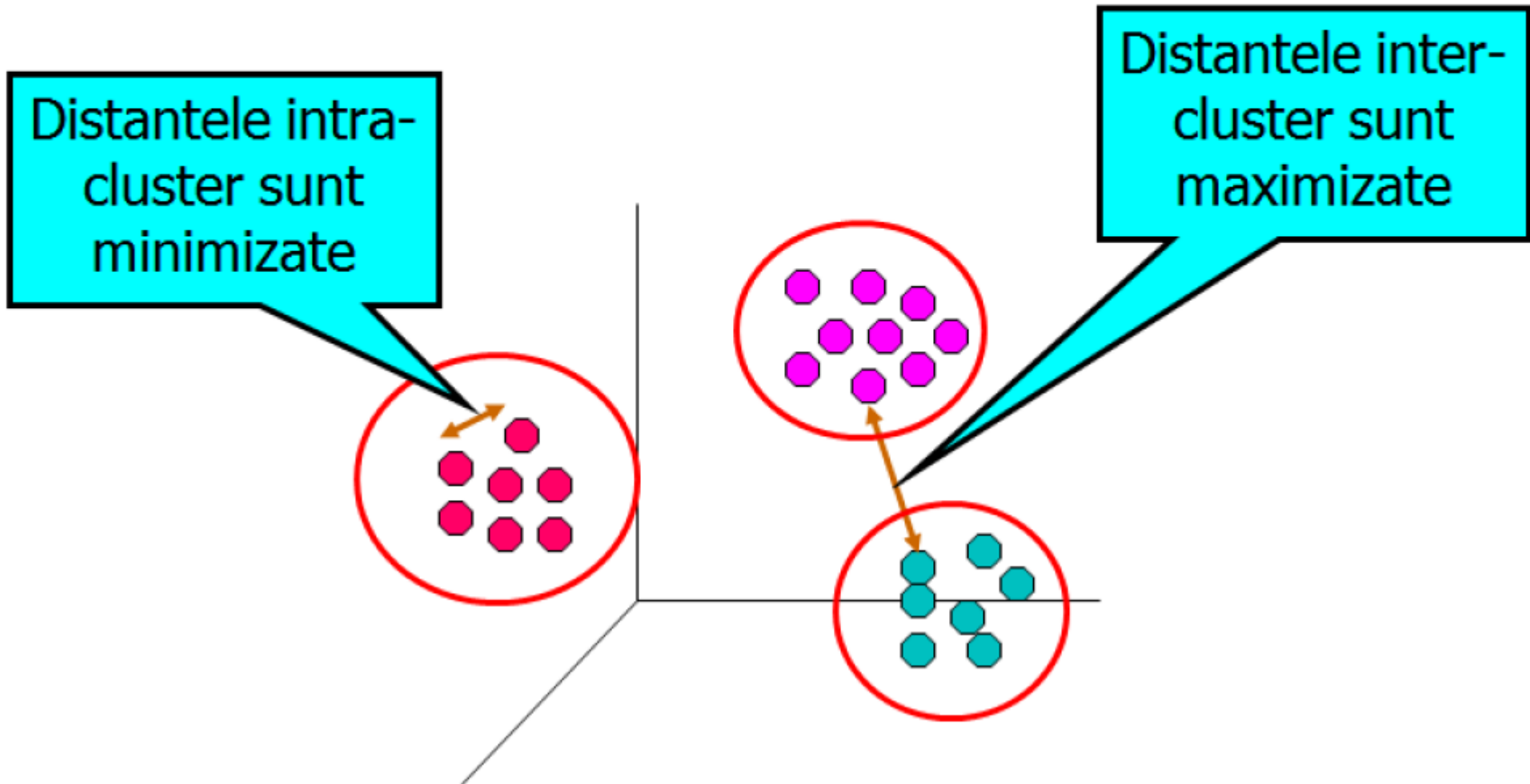
Gruparea (clusterizarea)

- Are ca scop găsirea unor grupuri astfel încât instanțele din același grup să fie mai asemănătoare între ele decât cu instanțele altor grupuri
- De obicei, este **nesupervizată**



Scop

- Analiza grupărilor (eng: cluster analysis) divide datele în grupuri, pe baza similarităților sau a relațiilor dintre ele
 - datele din același cluster sunt similare între ele
 - datele din alte cluster ar trebui să aibă similaritate mică cu cele din clusterul curent



Scopul analizei grupărilor

- Prin clustering se încearcă obținerea de grupări care sunt:
 - semnificative: clusterelor trebuie să surprindă natura structurală a datelor
 - utile: sumarizarea unui volum mare de date, furnizare de explicații sau ambele

Scop

Domenii de utilizare:

- științe sociale
- biologie
- statistică
- recunoaștere de șabloane (pattern recognition)
- regăsirea informației
- analiza imaginilor
- bioinformatică

Clustering pentru facilitarea înțelegerii

- Scop: obținerea de clase de elemente, utile în analizarea și înțelegerea lumii înconjurătoare.
- Exemple:
 - **biologie:** obținerea de taxonomii: regn, încregătură, clasă, ordin, familie, gen, specie;
 - **bioinformatică:** gruparea genelor și a proteinelor, presupuse a avea funcționalitate similară
 - **psihologie și medicină:** boli sau stări au diferențe minore; gruparea lor e folosită pentru descoperirea de subcategorii – e.g. diverse tipuri de depresii
 - **afaceri:** gruparea clienților în funcție de similaritățile de achiziție și realizarea de reclame particularizate pe grupuri de clienți

Clustering pentru scopuri utilitare

- Permite abstractizare față de datele care intră în cluster
- Alternativ: clusterelor se pot caracteriza deseori prin intermediul prototurilor
- Din acest punct de vedere clustering-ul este modalitatea de determinare a celor mai reprezentativi centri de cluster
 - **sumarizare:** în loc de a considera datele de plecare, se poate lucra cu centroizii clusterelor; rezultatele pot avea acuratețe comparabilă
 - **compresie:** pentru fiecare cluster se consideră centroidul lui; se obține o aproximare bazată pe substituirea datei inițiale cu centroidul clusterului de care aparține; se mai numește și **vector quantization**
 - **găsirea eficientă a celor mai apropiați vecini:** în loc să se calculeze distanțe între toate perechile de obiecte, se poate restricționa calculul la doar obiectele din același cluster, sau clusterelor apropiate

Ce nu este analiza clusterelor?

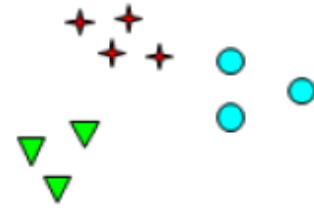
- Clasificare supervizată: clustering-ul este de fapt învățare nesupervizată
- Segmentare simplă a datelor, precum gruparea studenților după prima literă din nume
- Partiționarea de grafuri: există câteva elemente comune, dar partiționarea de grafuri se face după specificații externe mai complexe decât pentru clustering, sau subgrafurile obținute nu au o separare considerabilă

Ambiguitatea noțiunii de cluster

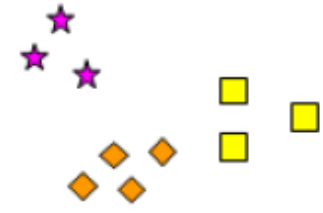
- În multe domenii noțiunea de cluster nu e clar definită



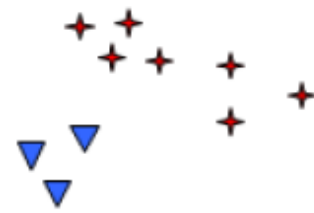
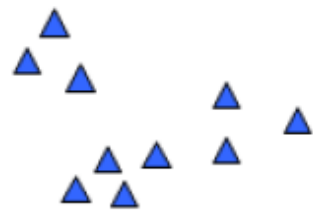
Cate cluster are?



Sase cluster



Doua cluster



Patru cluster



Figure: Moduri de grupare

Algoritmi de grupare

- ❑ Algoritmul k-medii (k-means)
- ❑ Aloritmul EM (Expectation-Maximization)
- ❑ Gruparea ierarhica

Algoritmi de grupare

- ❑ Algoritmul k-medii (k-means)
- ❑ Aloritmul EM (Expectation-Maximization)
- ❑ Gruparea ierarhica

Algoritmul k-medii (k-means)

- Algoritmul partiționează o mulțime de instanțe numerice în k grupuri (cluster)
- k reprezintă numărul de grupuri și este un parametru de intrare ales de utilizator

Algoritmul k-medii (k-means)

Modul de functionare

- Se inițializează aleatoriu cele k centre inițiale
- Se atribuie fiecărui centru instanțele cele mai apropiate de el

$$C(i) = \arg \min_{1 \leq k \leq K} \|x_i - m_k\|^2, \quad i = 1, \dots, N$$

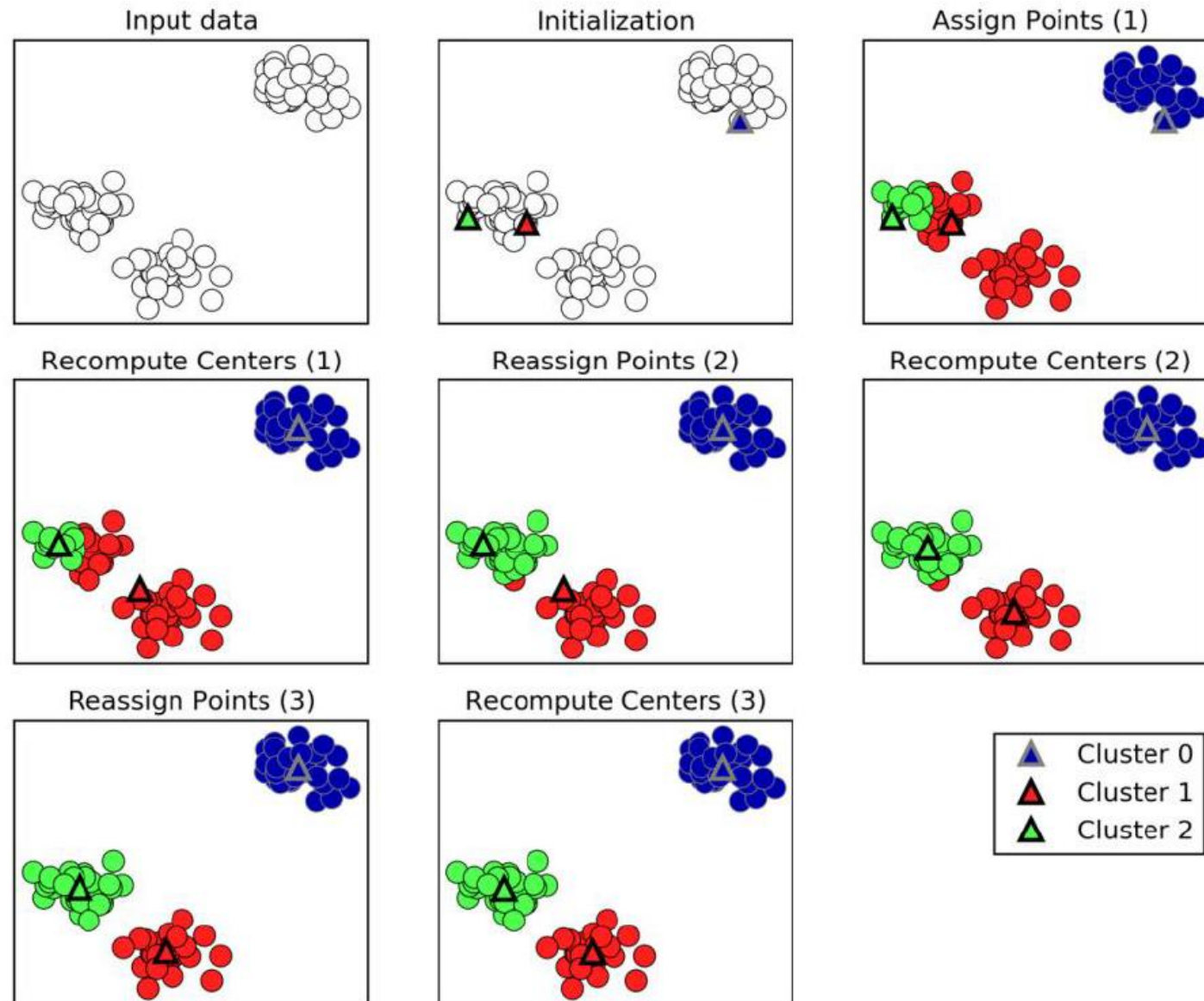
- Se calculează centrul de greutate (media aritmetică) a tuturor instanțelor atribuite unui centru și se actualizează poziția centrului grupului respective

$$m_k = \frac{\sum_{i:C(i)=k} x_i}{N_k}, \quad k = 1, \dots, K.$$

- Se repetă cei doi pași până când nu se mai modifică poziția niciunui centru

Algoritmul k-medii (k-means)

Exemplu



Algoritmul K-means

- Algoritm bazat pe prototipuri
- Prototipul este un centroid
- Tehnica se aplică pentru obiecte din spațiul continuu n -dimensional
- Cel mai popular algoritm de clustering
- Variantă înrudită: K-medoid — centroidul este unul din punctele din cluster

Pașii algoritmului K-means

- 1 Se aleg K centroizi inițiali; K e parametru de intrare
- 2 Fiecare punct este asignat celui mai apropiat centroid
- 3 Fiecare colecție de puncte asignată unui centroid este un cluster
- 4 Centroidul este modificat pe baza punctelor din clusterul pe care îl reprezintă
- 5 Se reia de la pasul 2 până când nu se mai schimbă poziția celor K centroizi

Algoritmul K-means: pseudocod

Algorithm 1 Basic K-means Algorithm.

- 1: Select K points as the initial centroids.
 - 2: **repeat**
 - 3: Form K clusters by assigning all points to the closest centroid.
 - 4: Recompute the centroid of each cluster.
 - 5: **until** The centroids don't change
-

Algoritmul K-means: observații

- Centrozii inițiali sunt cel mai des aleși aleator
- Centroidul este de regulă media punctelor din cluster
- Aproximarea dintre punct și centroid se consideră ca distanță sau similaritate
 - distanță Euclidiană, similaritatea cosinus, coeficient de corelație etc.
- K-Means va converge pentru măsurile de similaritate de mai sus
- Convergența este rapidă, apare după relativ puține operații
- De multe ori condiția de oprire se schimbă în: “până când relativ puține puncte își schimbă clusterul”
- Complexitatea este $O(m \cdot K \cdot l \cdot d)$, unde m = numărul de puncte, K = numărul de clustere, l = numărul de iterații, d = dimensiunea spațiului de intrare

Algoritmul K-means: exemplificare

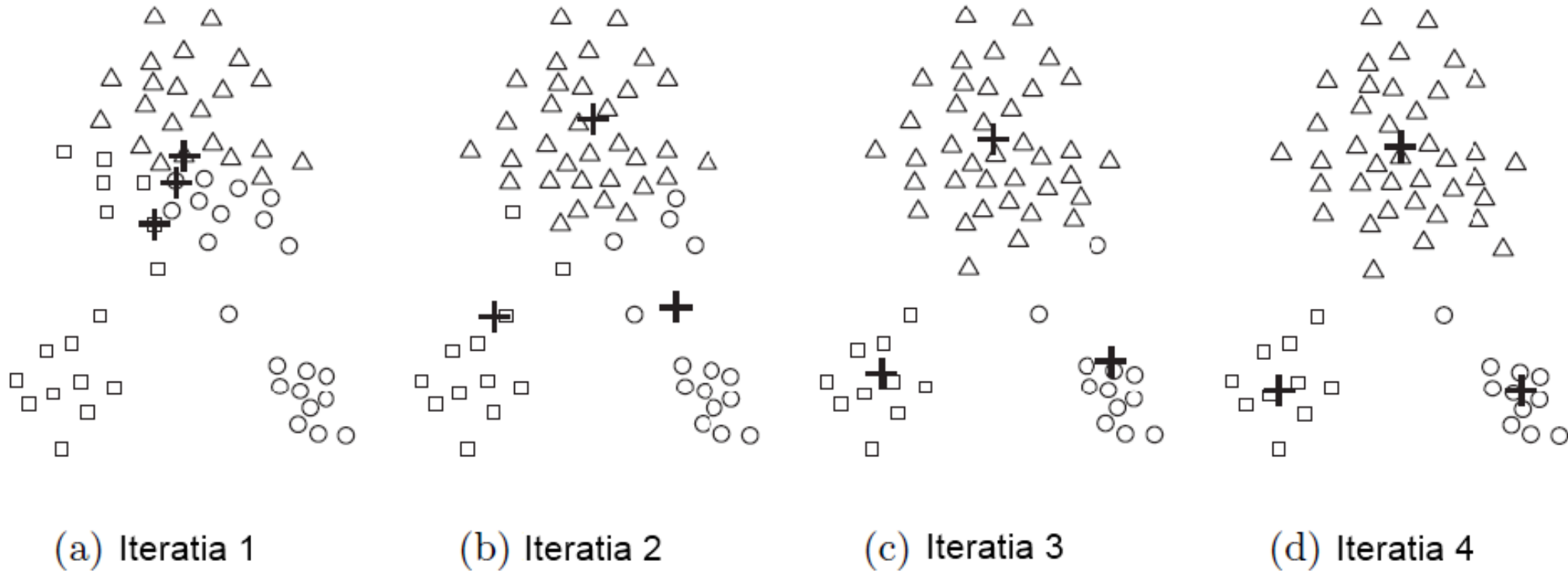


Figure: Exemplificare de iterații pentru K-means. Punctele marcate cu cruce sunt centroizii de la fiecare pas. La finalizare (pasul d) centroizii se stabilizează.

Algoritmul K-means: observații

- Pasul 4 din algoritm este:
Recompute the centroid of each cluster
deoarece centroidul se poate modifica
- Modificarea poziției centroidului duce la modificarea distanțelor de la punctele din cluster la centroid
- Problema de clustering poate fi privită ca una în care funcția obiectiv este de a minimiza suma pătratelor distanțelor de la puncte la centroizii clusterului din care fac parte
- Funcția obiectiv este utilizată pentru măsurarea calității clusteringului

Algoritmul K-means: clusteringul ca problemă de optimizare

- Funcția: **suma pătratelor erorilor**, orig: **sum of squared error (SSE)**, **scatter**
- Notatii:

Simbol	Descriere
\mathbf{x}	Un obiect
C_i	Al i -lea cluster
\mathbf{c}_i	Centroidul clusterului i
\mathbf{c}	Centroidul tuturor punctelor
m_i	Numărul de obiecte din al i -lea cluster
m	numărul de obiecte din setul de date
K	numărul de clustere presupuse

- Formula:

$$SSE = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} \text{dist}(\mathbf{c}_i, \mathbf{x})^2$$

Algoritmul K-means: clusteringul ca problemă de optimizare

- Dacă se formează clustere diferite atunci valoarea SSE poate diferi de la caz la caz
- Putem avea deci clustering optimal și suboptimal

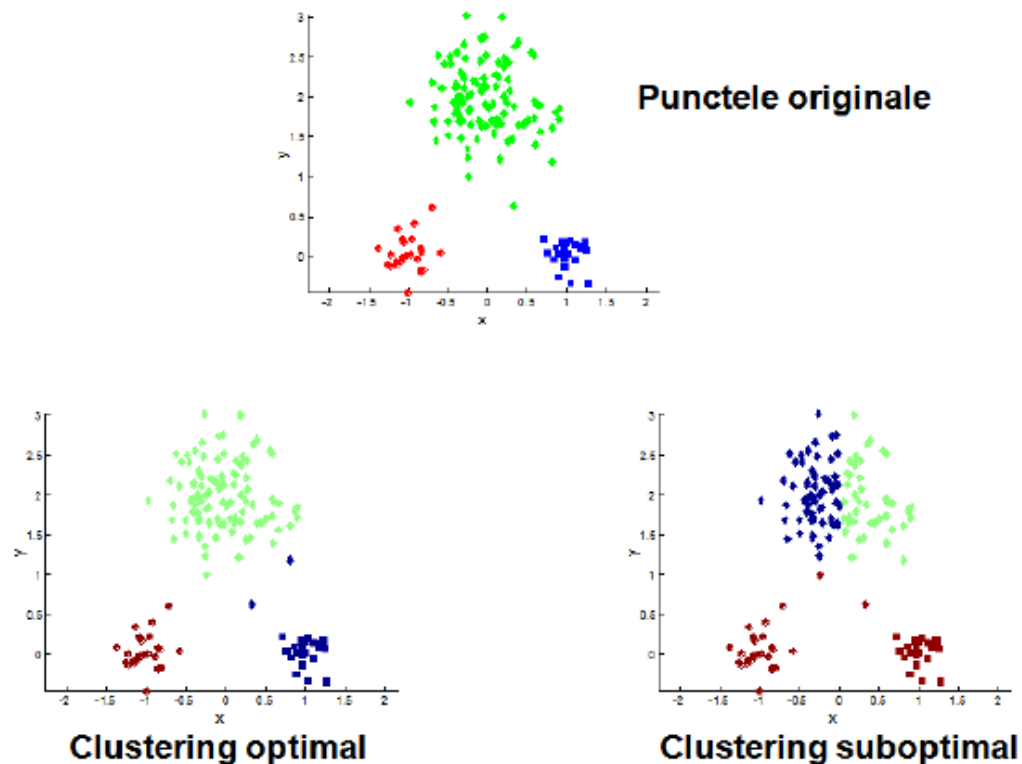


Figure: Clustering optimal vs. suboptimal, relativ la funcția SSE.

Algoritmul K-means: clusteringul ca problemă de optimizare

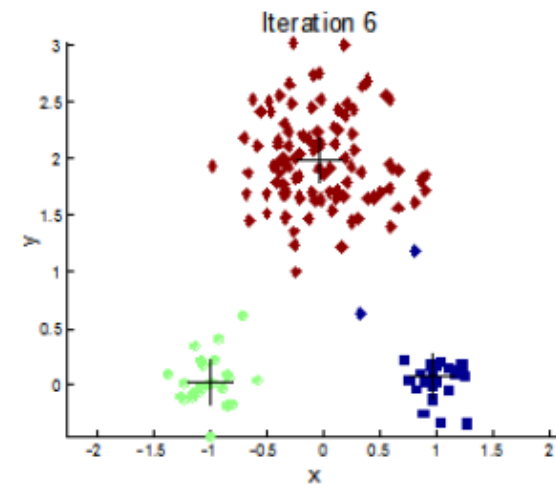
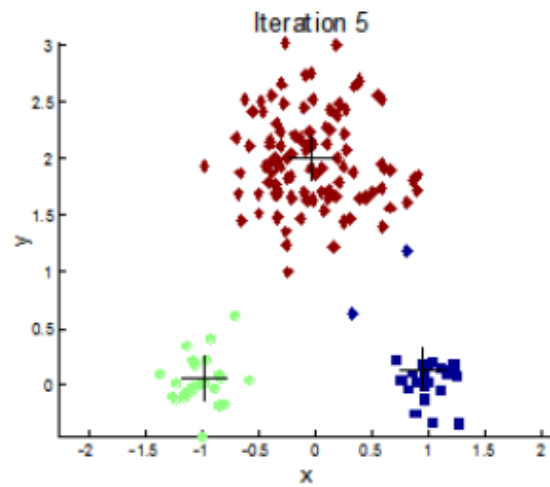
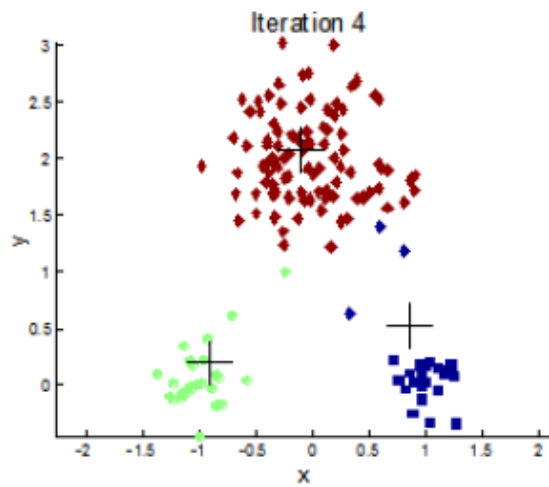
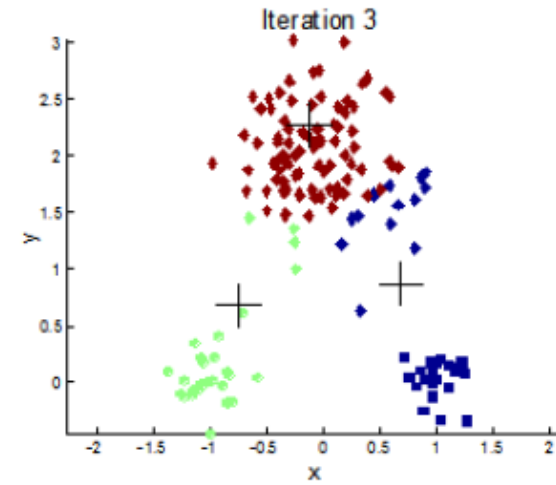
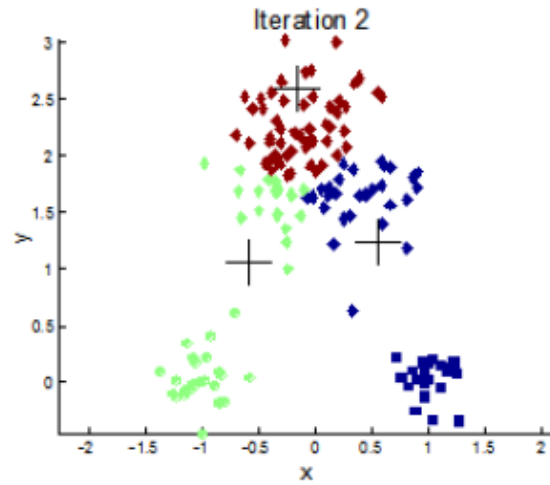
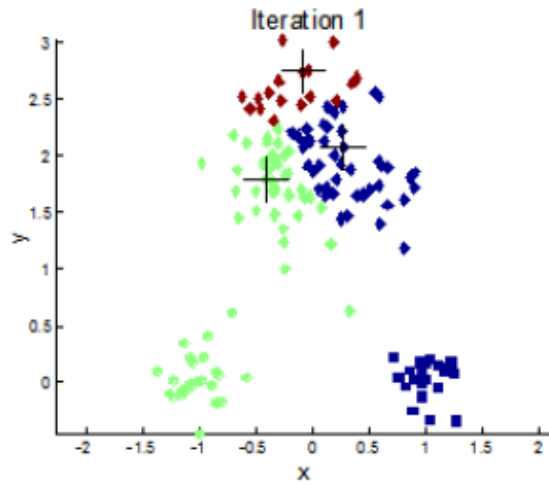
Algorithm 1 Basic K-means Algorithm.

- 1: Select K points as the initial centroids.
 - 2: **repeat**
 - 3: Form K clusters by assigning all points to the closest centroid.
 - 4: Recompute the centroid of each cluster.
 - 5: **until** The centroids don't change
-

- Pașii 3 și 4 din algoritm încearcă să minimizeze în mod direct SSE
- Optimul care se obține poate fi unul local: *se optimizează SSE pentru anumite alegeri ale locațiilor centroizilor, nu pentru toate variantele*

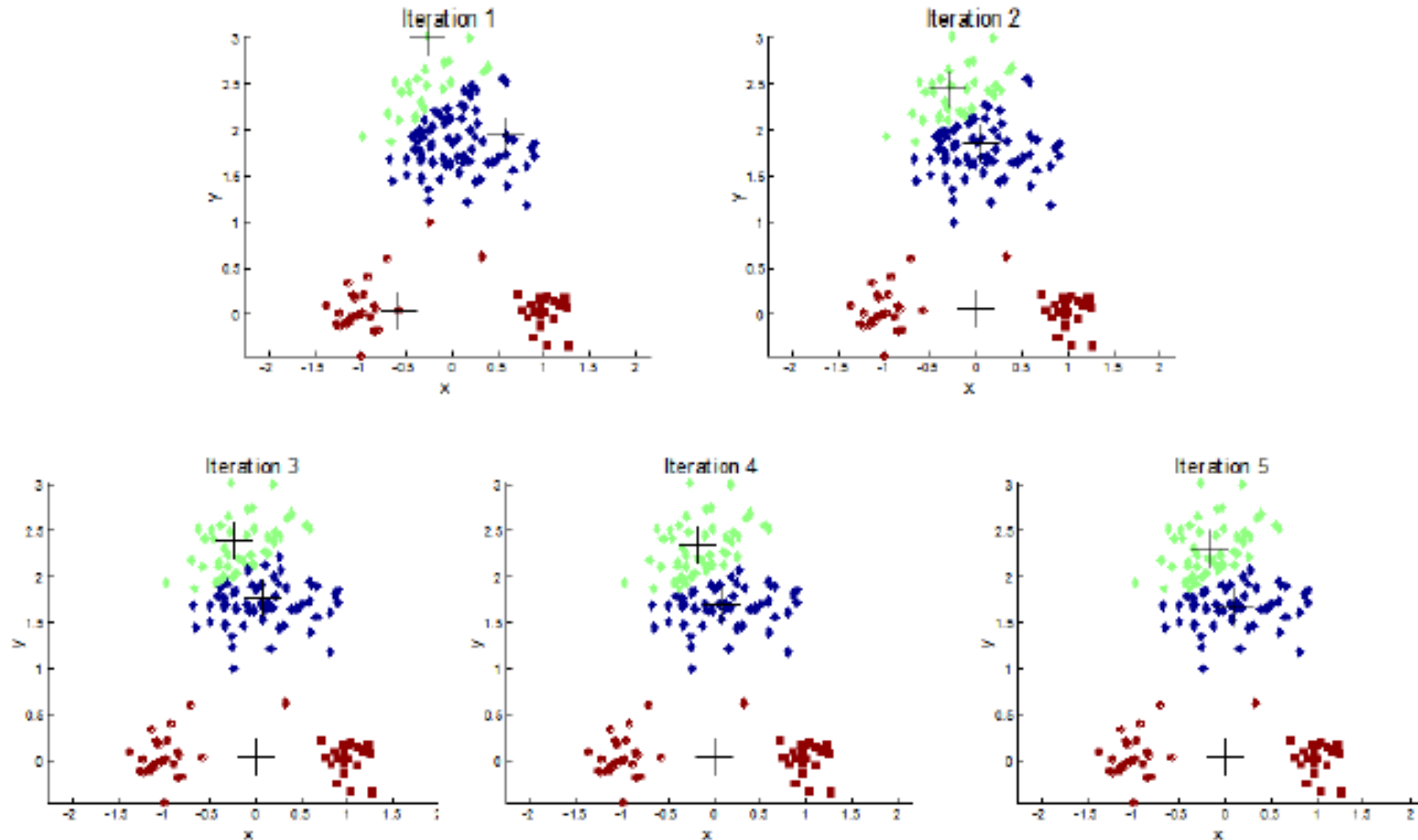
K-means: importanța alegerii centroizilor inițiali

- Centroizi inițiali aleși “inspirat”



K-means: importanța alegerii centroizilor inițiali

- Clustere inițiale alese nefavorabil



Algoritmul K-means: clusteringul ca problemă de optimizare

- K-means nu se reduce doar la date în spațiul Euclidian
- Pentru cazul în care datele se referă la documente: obiectivul este maximizarea similarității între documentele din același cluster = maximizarea **coeziunii** clusterului:

$$\text{Coeziunea totală} = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} \cos(\mathbf{x}, \mathbf{c}_i)$$

Algoritmul K-means: metode de reducere a impactului negativ de alegere nefavorabilă a centroizilor

- Se ia un eșantion de puncte și se aplică clustering ierarhic; centroizii rezultați sunt apoi folosiți ca și centroizi inițiali pentru K-means clustering
- K-means++ [1]
 - 1 Se alege uniform un centru din datele inițiale
 - 2 Pentru fiecare punct x se calculează distanța $D(x)$ de la x la cel mai apropiat centroid care a fost deja ales
 - 3 Alege un punct aleator, cu probabilitatea proporțională cu $D^2(x)$
 - 4 Repetă pașii 2 și 3 până când s-au ales K centroizi
- Bisecting K-means
- Postprocesare

Algoritmul K-means: problema clusterelor goale

- Algoritmul original poate duce la crearea de cluster goale
- Cluster gol: niciun punct nu este alocat aceluia cluster
- Dacă acest centroid este lăsat așa, SSE crește mai mult decât e cazul
- Soluții:
 - se renunță la acest centroid, se alege un altul (*e.g.* cel mai îndepărtat de oricare din ceilalți centroizi)
 - se alege un centroid din clusterul care contribuie cel mai mult la creșterea valorii SSE
- Dacă sunt mai multe cluster goale, se procedează în mod repetat ca mai sus

Algoritmul K-means: outliers

- Datele de tip outlier influențează calculul centroizilor; aceștia pot să nu mai fie reprezentativi
- De regulă se preferă eliminarea lor înainte de clustering
 - totuși: valorile outlier pot corespunde unor date de interes major
- Alternativ: se pot detecta și elimina outliers în faza de postprocesare:
 - se determină acele date care contribuie cel mai mult la creșterea SSE
 - se formează clustere mici constând în grupări de outliers

Algoritmul K-means: reducerea SSE prin postprocesare

- Strategii:
 - clusterul care are valoarea SSE cea mai mare poate fi împărțit în mai multe subcluster
 - se introduce un nou centroid: cel mai depărtat punct față de orice cluster
 - dispersarea unui cluster: se ignoră centroidul unui cluster, punctele din el sunt reasignate altor cluster; se alege clusterul care produce creșterea maximă de SSE
 - unirea a două cluster: clusterelor cu cel mai apropiat centroid se unesc și rezultă unul singur

Algoritmul K-means: modificarea incrementală a clusterelor

- In varianta originală poziția centroizilor este calculată doar după ce punctele sunt asignate celor mai apropiți centri
- Variantă: poziția centroidului să se calculeze după fiecare asignare a unui punct la acel centroid
- Avantaj: nu se ajunge la clusterare vide
- Se poate decide asignarea unui punct la un alt centroid doar dacă manevra duce la îmbunătățirea funcției obiectiv
- Dezavantaj: algoritmul devine dependent de ordinea de prezentare a datelor

Bisecting K-means

- Se bazează pe împărțirea succesivă a unui cluster în două subcluster
- Alegerea clusterului care se împarte se poate face după mai multe criterii: cel mai larg cluster, cel care contribuie cel mai mult la SSE, combinație de aceste două criterii
- Pas opțional: se aplică K-means cu centroizii obținuți prin bisecting K-means

Algorithm 8.2 Bisecting K-means algorithm.

- 1: Initialize the list of clusters to contain the cluster consisting of all points.
 - 2: **repeat**
 - 3: Remove a cluster from the list of clusters.
 - 4: {Perform several “trial” bisections of the chosen cluster.}
 - 5: **for** $i = 1$ to *number of trials* **do**
 - 6: Bisect the selected cluster using basic K-means.
 - 7: **end for**
 - 8: Select the two clusters from the bisection with the lowest total SSE.
 - 9: Add these two clusters to the list of clusters.
 - 10: **until** Until the list of clusters contains K clusters.
-

Bisecting K-means: exemplu

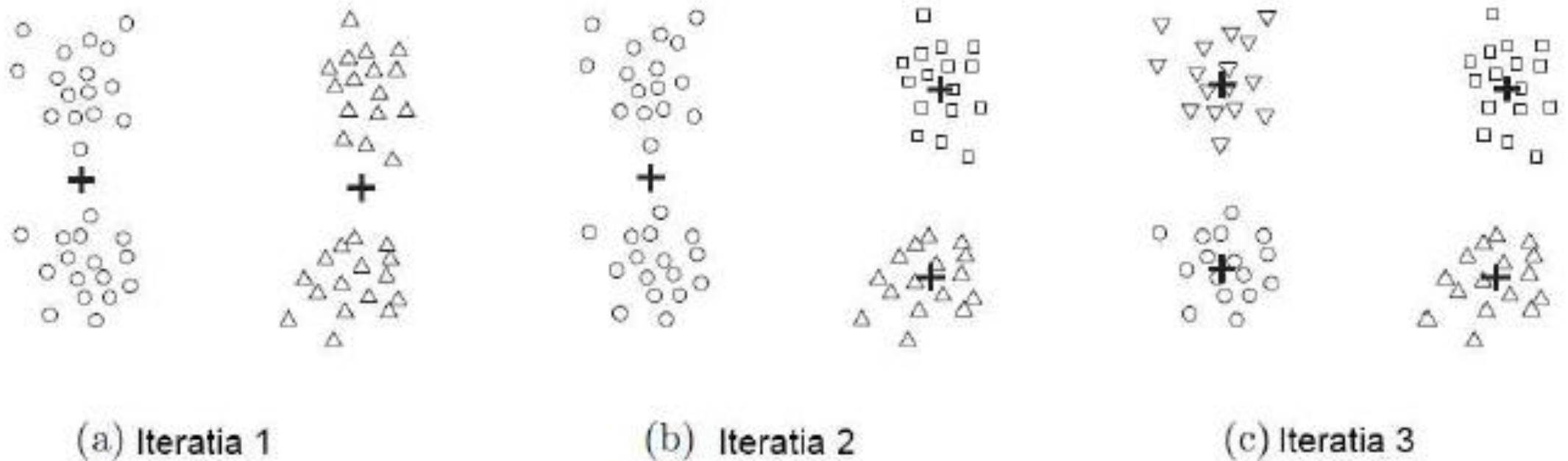
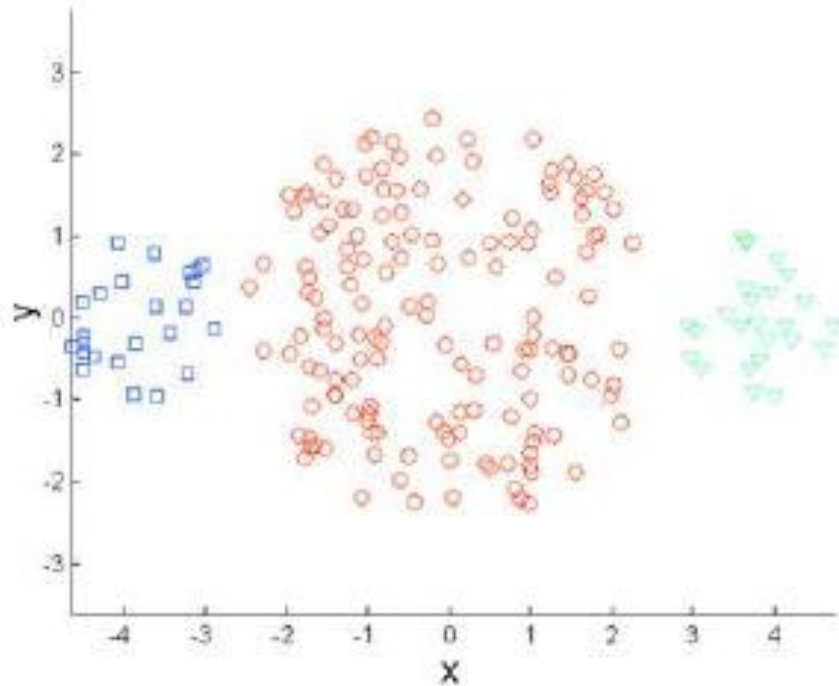


Figure: Exemplu de aplicare a algoritmului bisecting K-means.

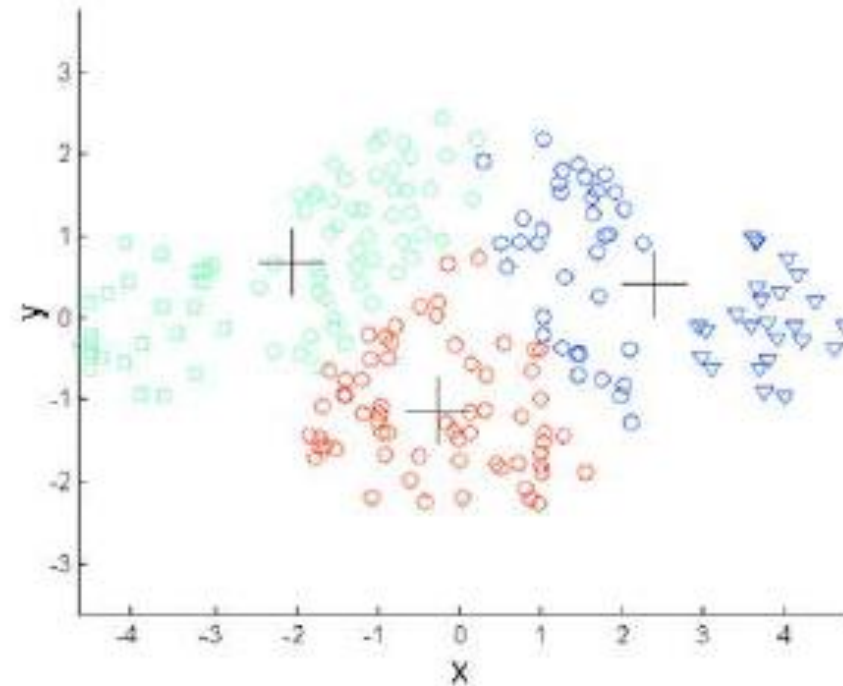
Limitări ale algoritmului K-means

- K-means intampină probleme cand clustererele sunt de diferite
 - dimensiuni
 - densități
 - forme, altceva decat cele circulare
- K-means are probleme cand datele prezintă outliers

Limitări ale K-means: cluster de dimensiuni diferite

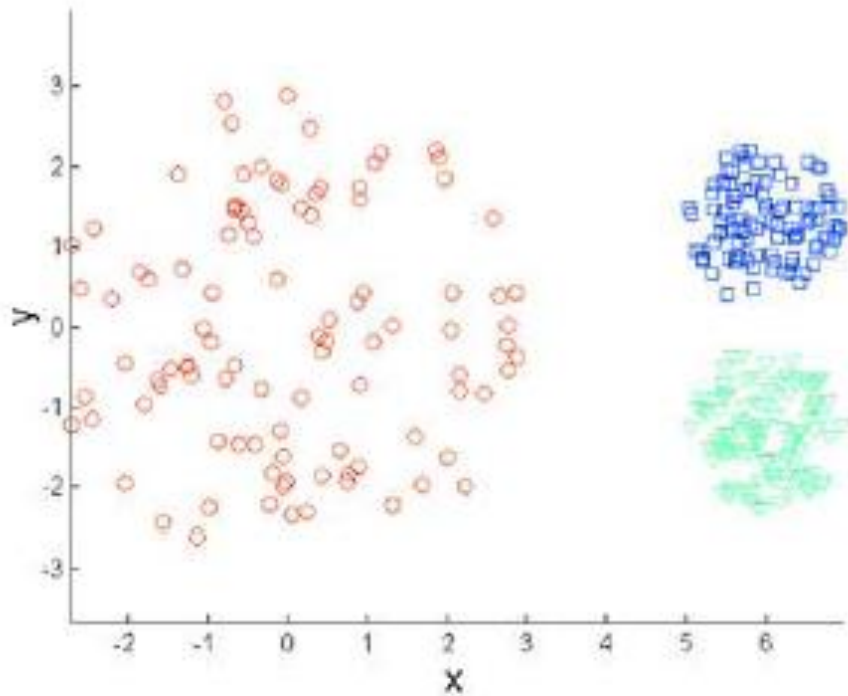


Punctele initiale

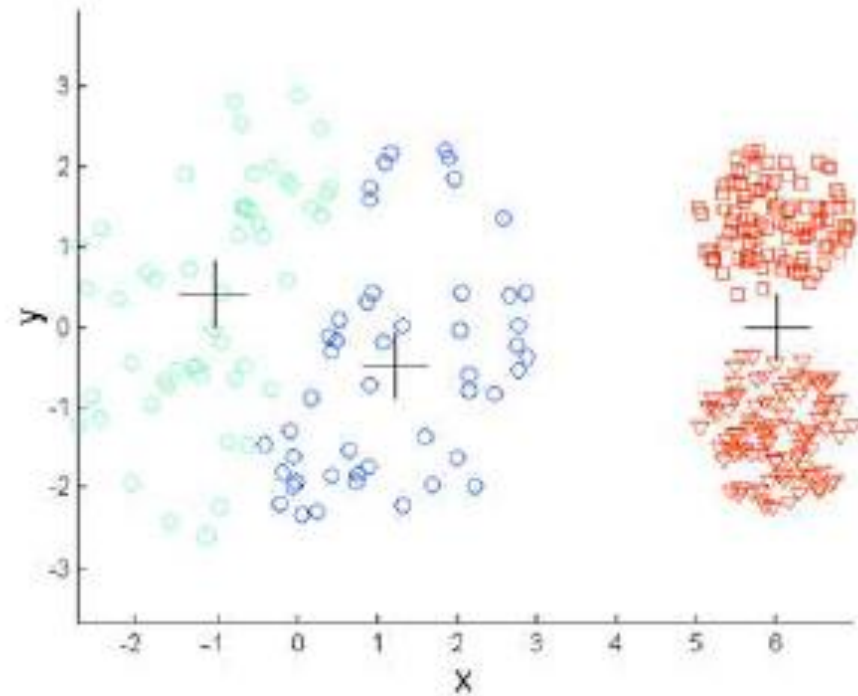


K-means (3 cluster)

Limitări ale K-means: clusterare de densități diferite

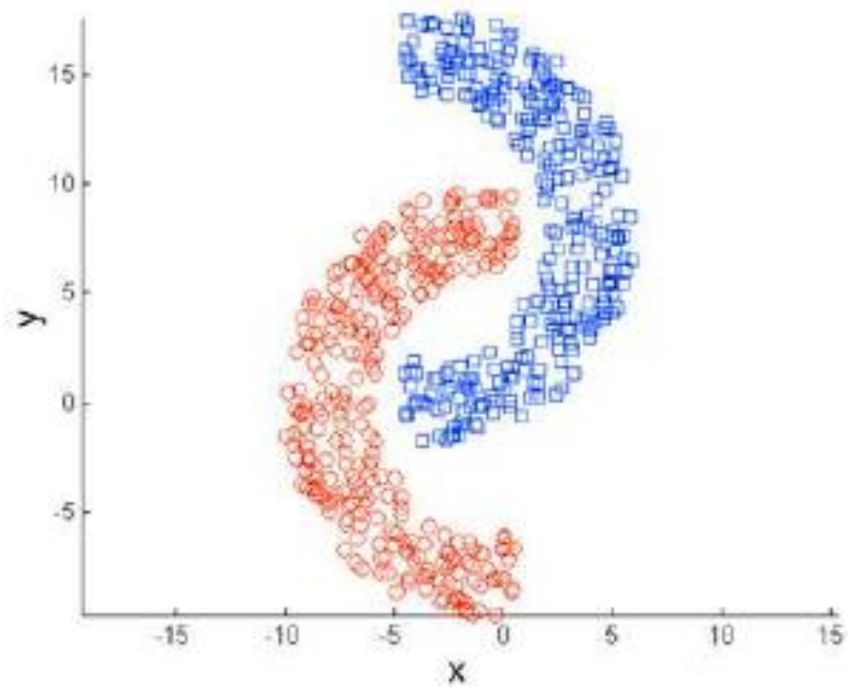


Punctele initiale

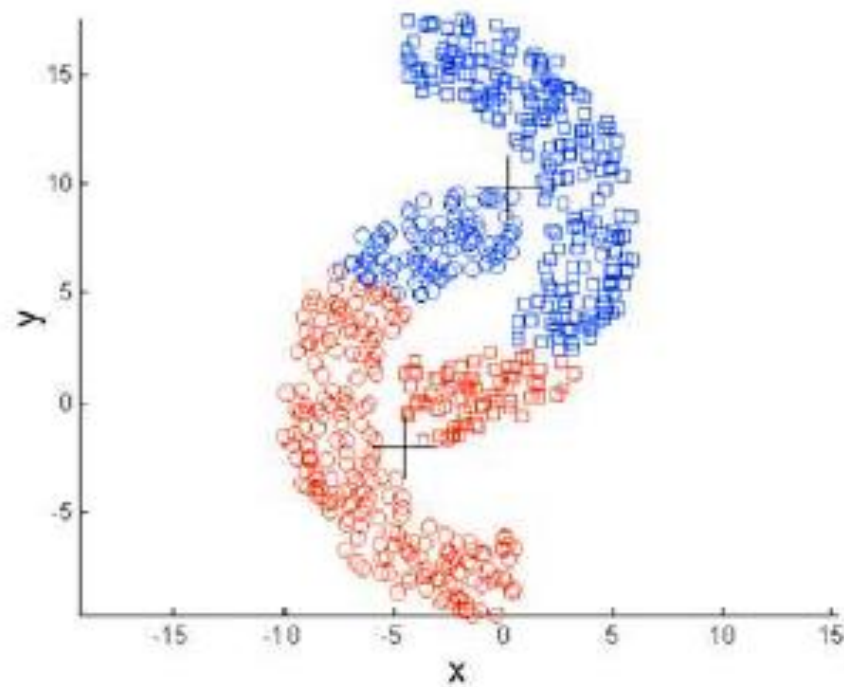


K-means (3 clusterare)

Limitări ale K-means: clusterare neglobulare



Punctele initiale



K-means (2 clusterare)

K-means: puncte tari și slabe

- Puncte tari:
 - simplu de implementat
 - eficient: se oprește după puține iterații
 - bisecting K-means, K-mean++: puțin sensibile la problema inițializării clusterelor
- Puncte slabe:
 - nu poate manipula date ce prezintă grupuri non-globulare, de dimensiuni sau densități diferite
 - probleme la datele care conțin outliers
 - algoritmul e restricționat la datele pentru care noțiunea de centroid are sens

Cuprins

- Generalități
- Algoritmul K-means
- Clustering ierarhic aglomerativ

Tipuri de clustering

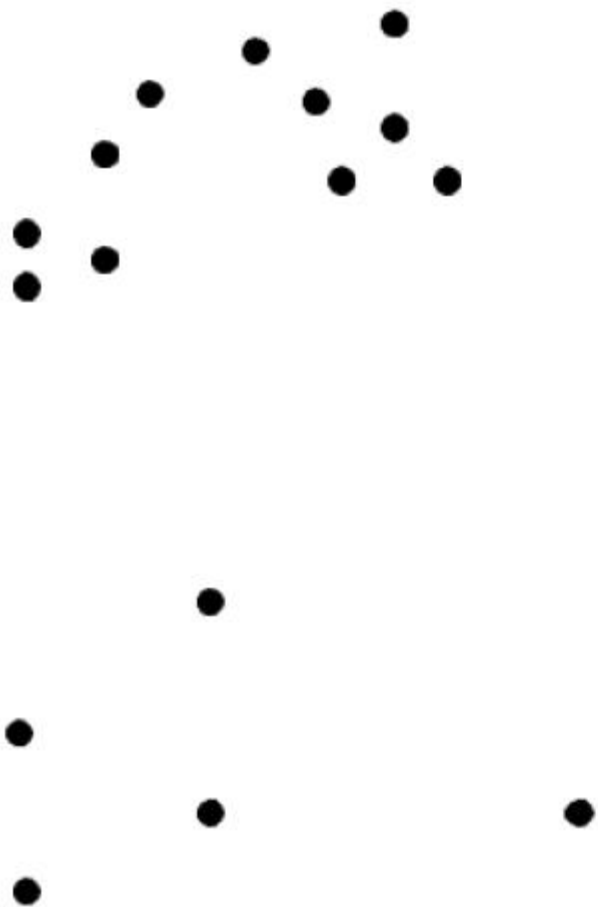
- Trei attribute ale unui proces de clustering:
 - ① ierarhic vs. partițional
 - ② exclusiv vs. cu suprapuneri vs. fuzzy
 - ③ complet vs. parțial

Clustering ierarhic vs. partițional

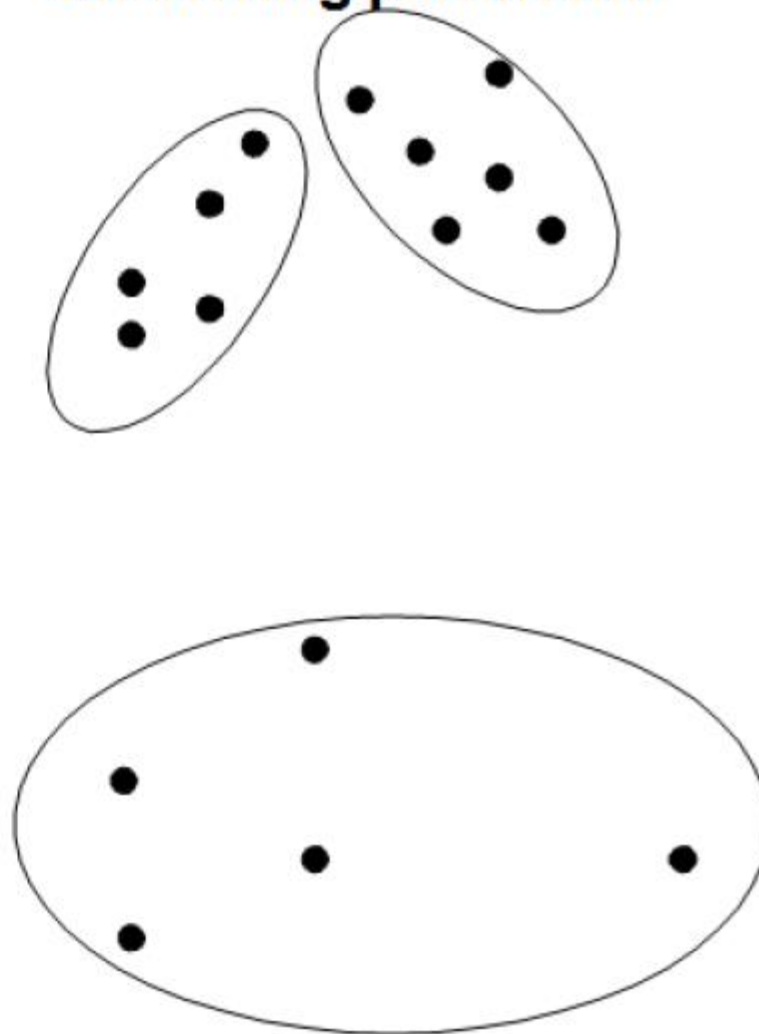
- Gruparea poate fi cu imbricări (ierarhică) sau nu (partițională)
- Clustering partițional: datele sunt împărțite în subseturi fără suprapuneri; fiecare dată face parte din exact un cluster
- Dacă permitem ca un cluster să aibă subclusterare \Rightarrow **clustering ierarhic**
 - fiecare nod intern este reuniunea clusterelor reprezentate de nodurile copil
 - se poate ajunge ca nodurile frunză să conțină exact o înregistrare
 - o partiționare ierarhică poate fi văzută ca o secvență de partiționări aplicată pornind de la mulțimea originală a datelor și grupând succesiv fiecare partiție în parte; partiționarea poate fi făcută până când se ajunge la noduri cu un singur element sau se poate face retezare

Clustering partiționat

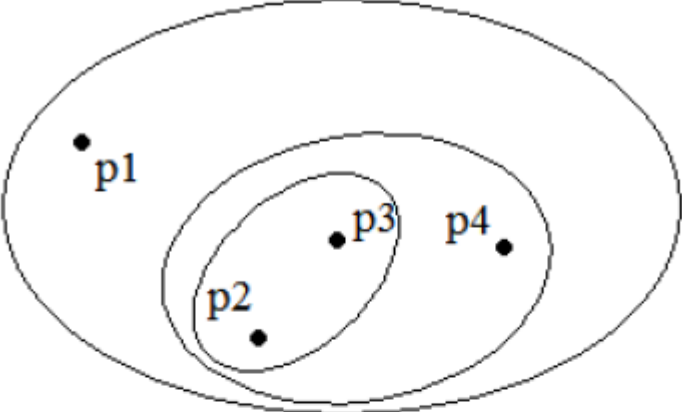
Setul de date



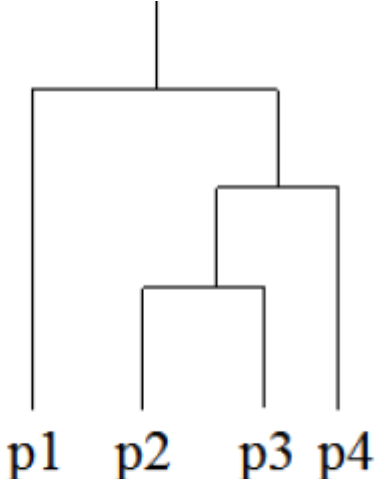
Clustering partiționat



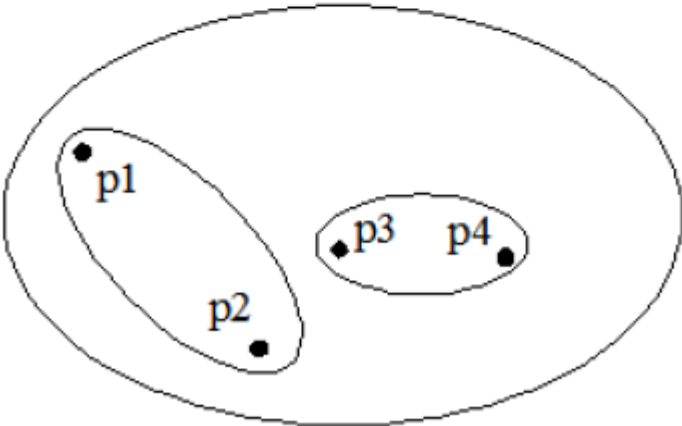
Clustering ierarhic



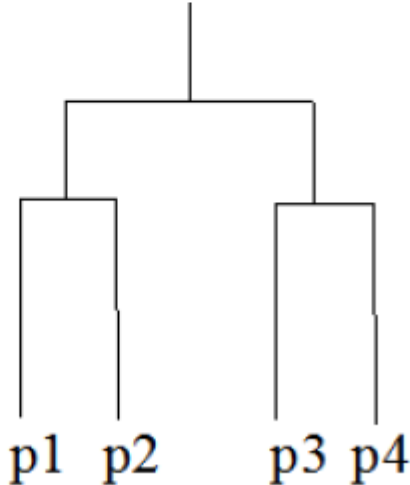
Clustering ierarhic traditional



Dendrograma traditionala



Clustering ierarhic netraditional



Dendrograma netraditionala

Clustering exclusiv vs. cu suprapuneri vs. fuzzy

- Exclusiv vs non-exclusiv:
 - exclusiv: un punct aparține unui singur cluster
 - non-exclusiv (eng: overlapping): un punct poate fi asociat mai multor cluster
 - utilitatea clusterelor ne-exclusive: pot reprezenta puncte ce aparțin simultan mai multor clase sau puncte apropiate de zona de separare
- Fuzzy clustering: fiecare obiect are o măsură fuzzy de apartenență la fiecare cluster
 - clusterelor devin mulțimi fuzzy
 - clusterelor fuzzy pot fi convertite la unele de tip exclusiv prin alegerea pentru fiecare dată a celui cluster pentru care măsura fuzzy este mai mare

Clustering parțial vs. complet

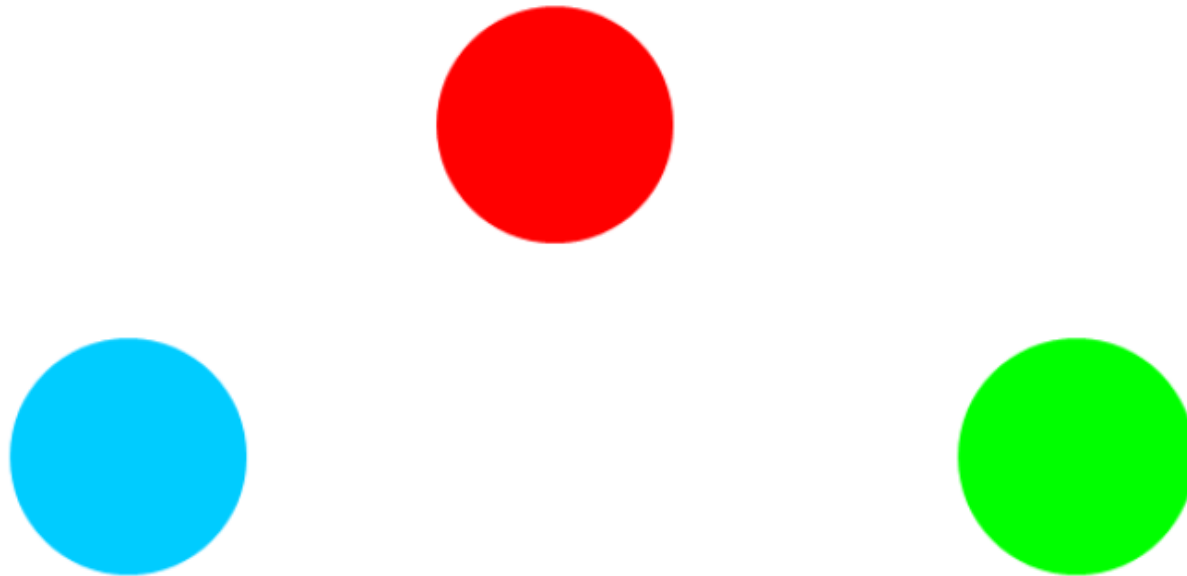
- Clustering complet: fiecare dată este asociată unui cluster
- Clustering parțial: pentru un obiect se poate să nu avem asignare la un cluster

Tipuri de clusterare

- Bine separate
- Bazate pe prototipuri
- Bazate pe grafuri
- Bazate pe densitate
- Conceptuale

Clustere bine separate

- În această viziune: un cluster este un set de puncte astfel încât orice punct din cluster este mai aproape (sau mai similar) față de orice punct din acel cluster decât de puncte din afara lui
- Definiția e respectată doar când datele sunt grupate grupate natural în clustere care sunt depărtate unele de altele



3 clustere bine separate

Clustere bazate pe prototipuri

- În această viziune: un cluster e o mulțime cu proprietatea că un obiect din cluster este mai apropiat/similar de/cu prototipul clusterului decât cu alte prototipuri
- Ca prototip: centroid sau medoid
 - medoid — punct din setul de date inițial care e cel mai reprezentativ pentru cluster
 - centroid — e.g. centru de greutate; valoarea lui poate să nu coincidă cu niciunul din punctele din centroid



4 clustere bazate pe centri

Clustere bazate pe grafuri

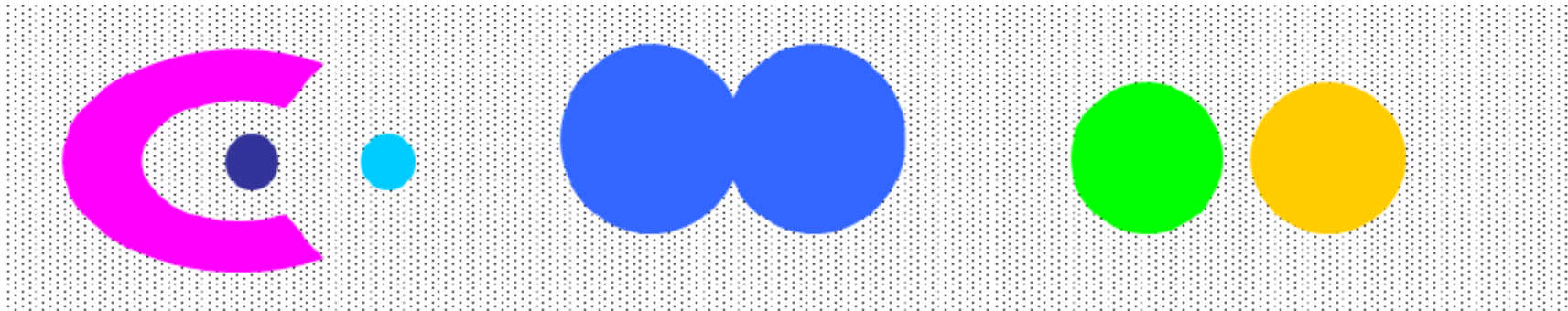
- Util pentru date reprezentate ca graf; un cluster este un grup de obiecte conectate
- Exemplu: clustere bazate pe contiguitate: două date sunt conectate doar dacă se află la o distanță mai mică decât un prag impus
- Posibile probleme: zgomotul poate crește clusterelor în mod artificial - vezi clusterul cu albastru deschis



8 clustere contigue

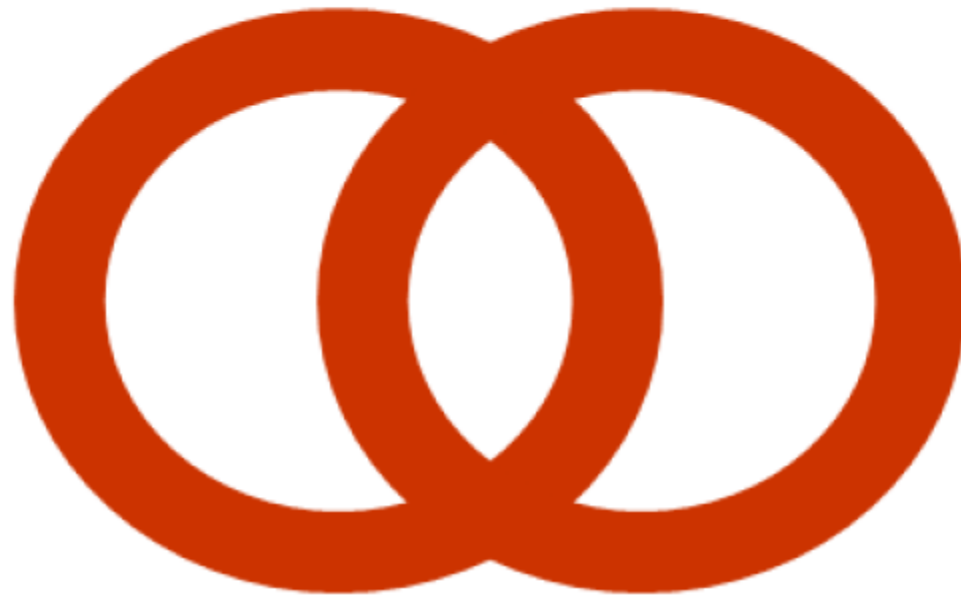
Clusterare bazate pe densitate

- În această viziune: un cluster este o zonă densă înconjurată de zonă de densitate mică
- Zonele de densitate mică separă pe cele cu densitate mare
- Utilizate când clusterurile sunt neregulate sau se înconjoară unele pe altele sau când zgomotul și datele outlier sunt prezente



Clustere conceptuale

- Clusterul este văzut ca un set de obiecte care partajează o proprietate ce reiese din mulțimea de puncte din cluster
- Un algoritm de clustering de acest tip ar trebui să detecteze concepte



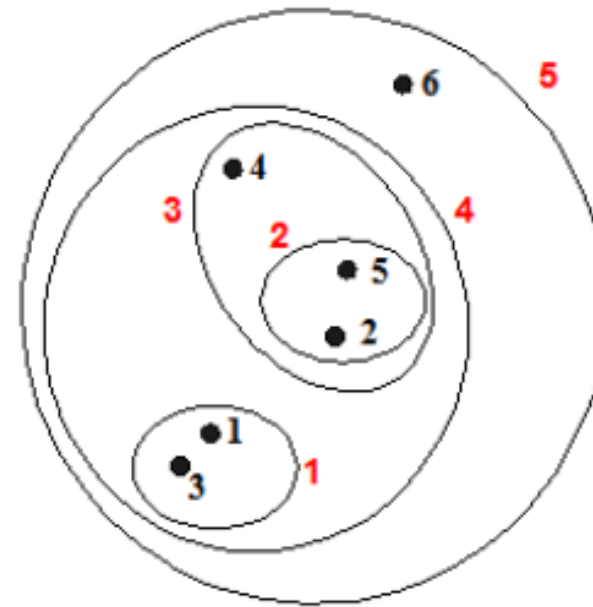
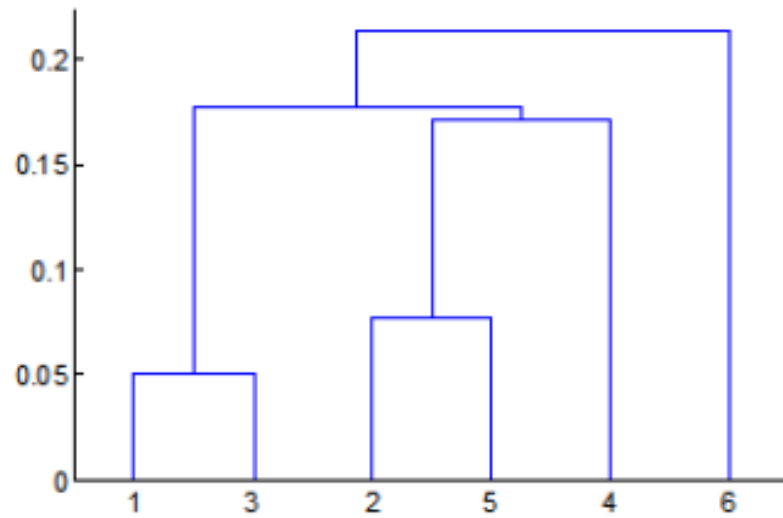
2 clustere sub forma de cercuri

Cuprins

- Generalități
- Algoritmul K-means

Clustering ierarhic

- Produce un set de clustere imbricate organizate ca un arbore
- Clusterelor pot fi vizualizate sub forma unei dendrograme



Clustering ierarhic

- Nu trebuie făcute presupuneri asupra numărului de clustere
 - se poate obține orice număr de clustere prin retezarea dendrogramei la nivelul dorit
- Pot corespunde unei taxonomii
 - exemplu: taxonomia din științele biologice

Clustering ierarhic

- Două tipuri de clustering ierarhic
- **Aglomerativ:**
 - se pornește cu fiecare punct ca și cluster
 - la fiecare pas se unește cea mai apropiată pereche de clustere până răman k clustere
- **Diviziv:**
 - se pleacă cu un cluster ce conține toate punctele
 - la fiecare pas se divizează un cluster
 - se continuă procesul până cand fiecare cluster conține un punct, sau se face rețezare când se ajunge la k clustere

Clustering ierarhic aglomerativ

- Tehnica este populară și cu vechime

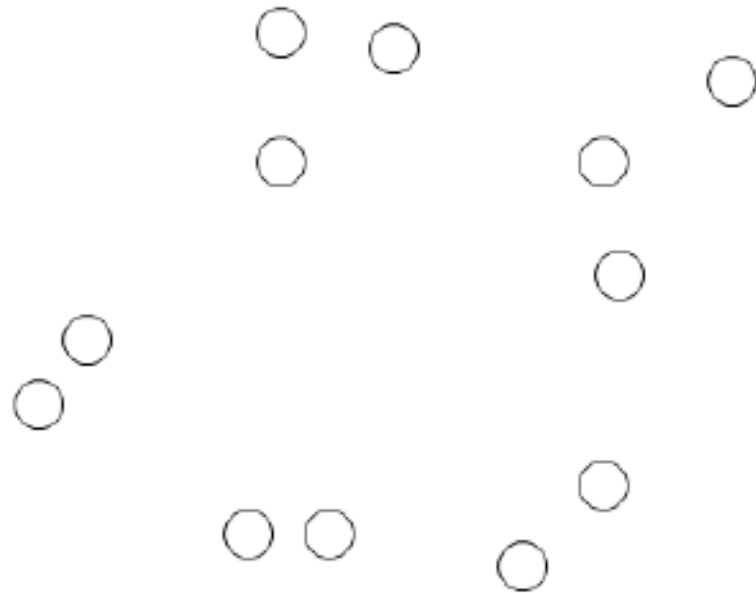
Algorithm 8.3 Basic agglomerative hierarchical clustering algorithm.

- 1: Compute the proximity matrix, if necessary.
 - 2: **repeat**
 - 3: Merge the closest two clusters.
 - 4: Update the proximity matrix to reflect the proximity between the new cluster and the original clusters.
 - 5: **until** Only one cluster remains.
-

- Operația esențială este calculul proximității a două clustere
- Diferențele între algoritmi ierarhici aglomerativi sunt date de modul de calcul al proximității

Clustering ierarhic aglomerativ: pornirea procesului

- Se pornește cu clustere formate din câte un punct



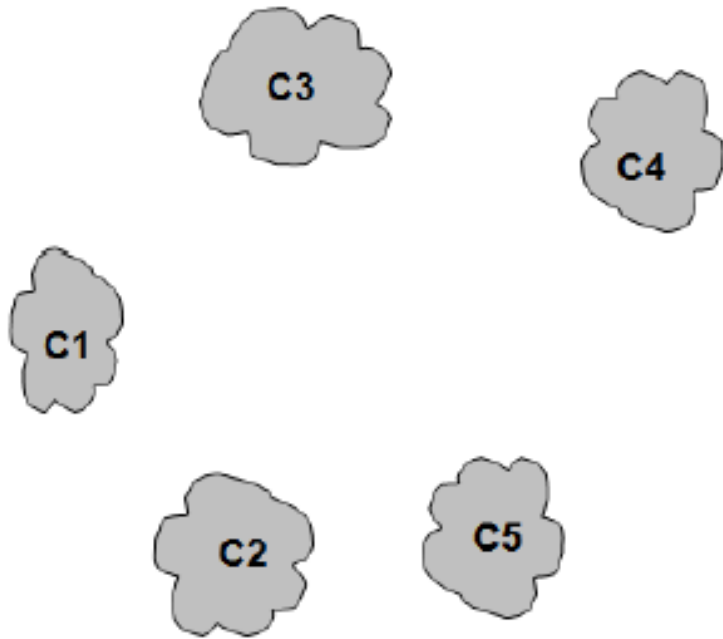
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Matrice de proximitate



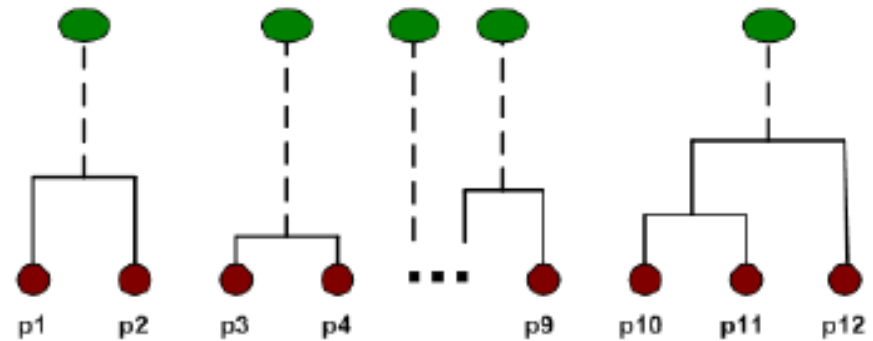
Clustering ierarhic aglomerativ: starea intermediară

- După câțiva pași avem câteva clustere



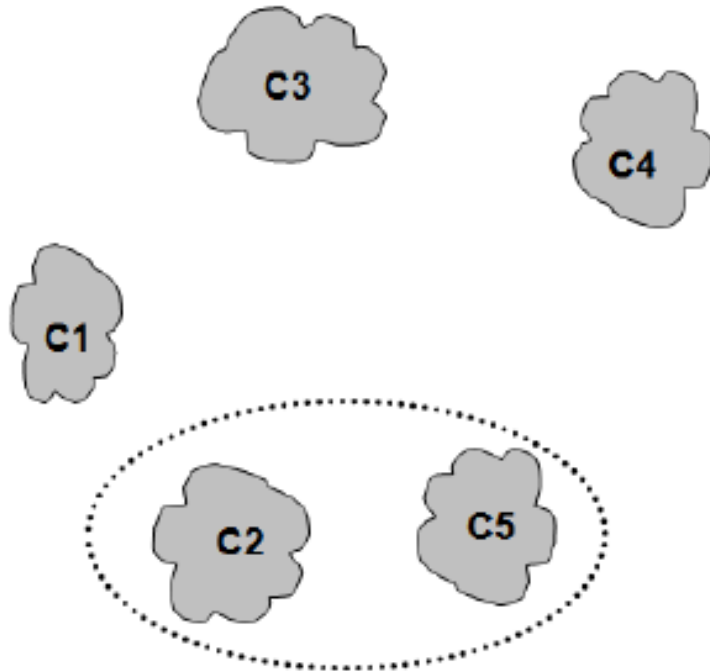
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Matrice de proximitate



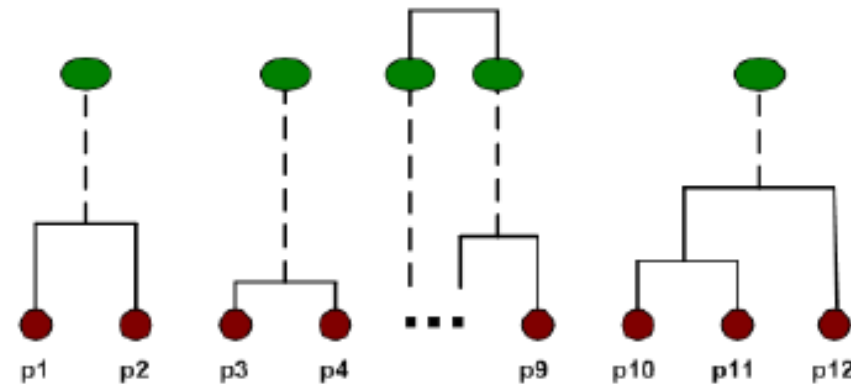
Clustering ierarhic aglomerativ: unirea de 2 cluster

- Unim cele mai apropiate două cluster $C2$ și $C5$ și modificăm matricea de proximitate



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Matrice de proximitate



Materiale de citit

- **Capitolul 10** din cartea: **An Introduction to Statistical Learning with Applications in R**. Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani. Springer-Verlag, 2013. (carte disponibilă gratuit online aici: <http://www-bcf.usc.edu/~gareth/ISL/>)

- Informațiile prezentate au fost colectate din diferite surse de pe internet, precum și din slide-urile de Data mining ale prof. Lucian Sasu, Univ. Transilvania din Brasov, disponibile pe internet.