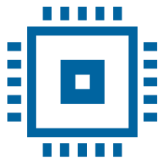


Analiza și procesarea datelor prin tehnici de Învățare Automată

5. Tehnici de învățare nesupervizată: Clustering



Universitatea
Transilvania
din Brașov

FACULTATEA DE INGINERIE ELECTRICĂ
ȘI ȘTIINȚA CALCULATOARELOR

Șef Lucrări Dr. ing. Horia Modran

Contact: horia.modran@unitbv.ro / modranhoria@gmail.com

Tel: 0770171577

2024 - 2025



Domenii de învățare automată

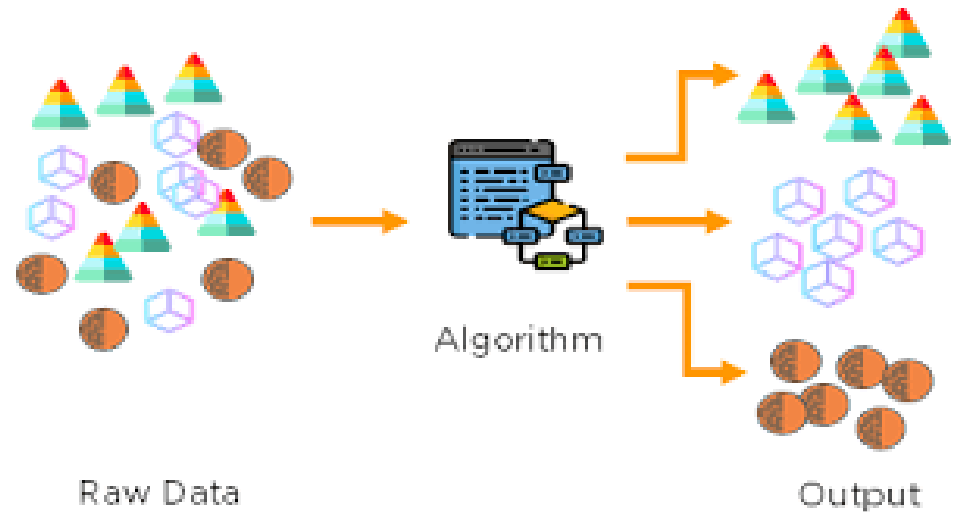
- **Învățare supervizată:** sunt furnizate datele și etichetele corespunzătoare
- **Învățare nesupervizată:** sunt furnizate doar date (nu sunt furnizate etichete)
- **Învățare semisupervizată:** sunt prezente unele etichete (dacă nu toate)
- **Învățarea prin întărire:** un agent care interacționează cu lumea face observații, întreprinde acțiuni și este recompensat sau pedepsit; acesta ar trebui să învețe să aleagă acțiunile astfel încât să obțină multe recompense





Învăţare nesupervizată

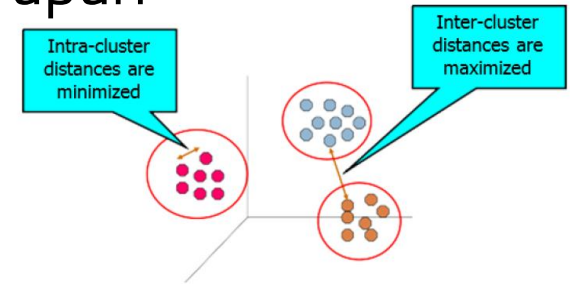
- datele nu au atribut ţintă
- dorim să explorăm datele pentru a găsi anumite structuri (şabloane) în ele
- putem efectua regresie/clasificare aici? **NU**
- metode: clustering





Clustering

- găsirea unor grupuri (clustere) de obiecte astfel încât obiectele dintr-un grup să fie similare (sau înrudite) între ele şi diferite (sau fără legătură) faţă obiectele din alte grupuri
- Analiza cluster nu este:
 - clasificare supervizată
 - segmentarea simplă (de exemplu, împărţirea elevilor în diferite grupuri de înregistrare în ordine alfabetică)
 - rezultatele unei interogări
 - grupările sunt un rezultat al unei specificaţii externe





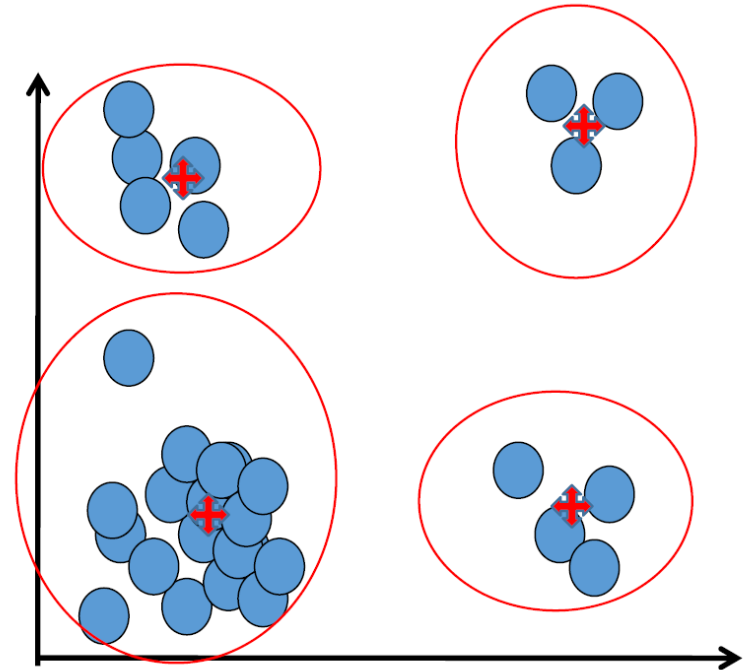
Clustering

- Clustering - tehnică de găsim a grupurilor similare în date, numite **cluster**
- acesta grupează instanţele de date care sunt similare (apropiate) între ele într-un singur cluster şi instanţele de date care sunt foarte diferite (îndepărtate) în cluster diferite.
- gruparea este adesea numită sarcină de învăţare nesupervizată, deoarece nu sunt date valori de clasă care să denotă o grupare a priori a instanţelor de date (cum este în învăţarea supervizată)
- clustering \sim învăţare nesupravegheată



Cluster

- un cluster este reprezentat de un singur punct, cunoscut sub numele de **centroid** al clusterului
- **centroidul** este calculat ca medie a tuturor punctelor de date dintr-un cluster $C_j = \sum x_i$
- limita clusterului este stabilită de cel mai îndepărtat punct de date din cluster (***borden***)





Aplicații

- Exemplul 1: grupează persoane de dimensiuni similare pentru a selecta dimensiunea tricourilor - "mici", "medii" și "mari"
 - personalizate pentru fiecare persoană: prea scumpe
 - one-size-fits-all: nu se potrivește tuturor
- Exemplul 2: în marketing, segmentarea clienților în funcție de asemănările lor
 - pentru a realiza marketing țintit
- Exemplul 3: Având în vedere o colecție de documente text, dorim să le organizăm în funcție de similitudinile conținutului lor pentru a produce o ierarhie tematică



Aspecte ale clustering

- Un algoritm de grupare a datelor
 - clustering partiţional
 - clustering ierarhic
 - clustering K-means
- o funcţie de distanţă (similaritate sau disimilaritate)
- calitatea grupării
 - distanţa între clustere \Rightarrow maximizată
 - distanţa intra-cluster \Rightarrow minimizată
- calitatea rezultatului tehnicii de clusering depinde de algoritm, de funcţia de distanţă şi de aplicaţie



Tipuri de clustering

- **Clustering:** sarcina de a grupa un set de puncte de date astfel încât punctele de date din acelaşi grup să fie mai asemănătoare între ele decât punctele de date din alt grup (grupul este cunoscut sub numele de cluster)
- **Tipuri:**
 - 1. Clustering ierarhic: Clustering aglomerativ, clustering diviziv
 - 2. Clustering exclusiv: K-means
 - 3. Clustering bazat pe densitate: DBSCAN



Tipuri de clustering

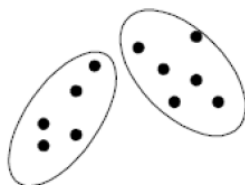
- distincţie importantă între clustere **ierarhice** şi **partiţionale**
- Clustering partiţional
 - împărţire a obiectelor de date în subseturi (grupuri) care nu se suprapun, astfel încât fiecare obiect de date să se afle în exact un subset
- Gruparea ierarhică
 - un set de clustere imbricate organizate sub forma unui arbore ierarhic



Partiţional vs. ierarhic

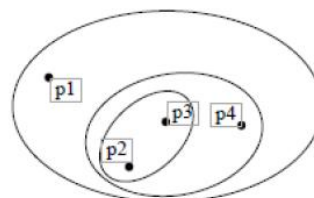
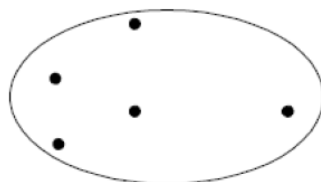


Original Points

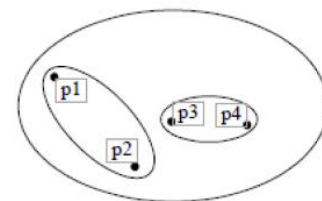


A Partitional Clustering

Clustering partiţional

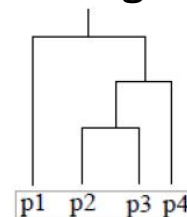


Traditional Hierarchical Clustering

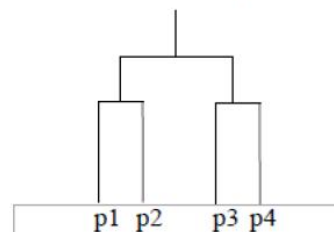


Non-traditional Hierarchical Clustering

Clustering ierarhic



Traditional Dendrogram



Non-traditional Dendrogram



Similaritate și disimilaritate

■ Similaritate

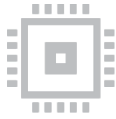
- măsură numerică a asemănării dintre două obiecte
- este mai mare atunci când obiectele sunt mai asemănătoare
- se încadrează adesea în intervalul $[0,1]$

■ Disimilaritate

- măsură numerică a diferenței dintre două obiecte de date
- mai mică atunci când obiectele sunt mai asemănătoare
- disimilaritatea minimă este adesea 0

■ Proximitatea se referă la o similitudine sau disimilaritate





Similaritate şi disimilaritate

- p şi q sunt valorile atributelor pentru două obiecte de date

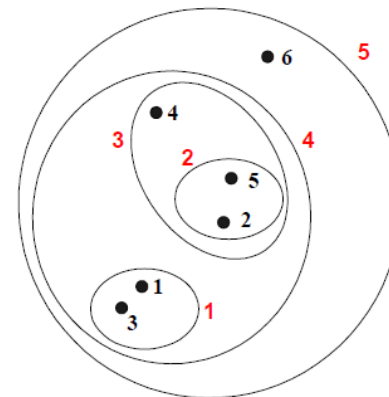
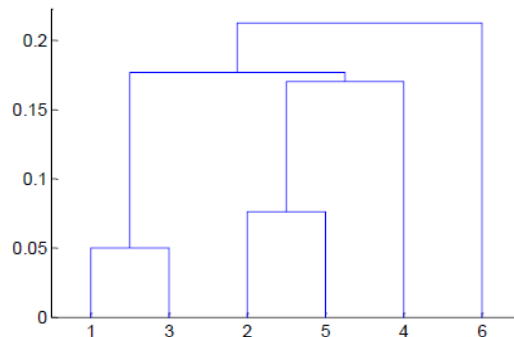
Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n-1$, where n is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval or Ratio	$d = p - q $	$s = -d, s = \frac{1}{1+d}$ or $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

Similarity and dissimilarity for simple attributes



Clustering ierarhic

- produce un set de clustere imbricate organizate ca un arbore ierarhic
- poate fi vizualizată ca o dendrogramă
 - o diagramă arborescentă care înregistrează secvenţele de fuziuni sau divizări





Tipuri

- Două tipuri principale de clustering ierarhic
 - Aglomerativ:
 - la început fiecare punct va reprezenta un cluter
 - la fiecare pas, se fuzionează cea mai apropiată pereche de clustere până când rămâne un singur cluster
 - Diviziv:
 - se începe cu un singur cluster ce conţine toate punctele
 - la fiecare etapă, se crează un nou cluster până când fiecare cluster conţine un punct (sau există k clustere)





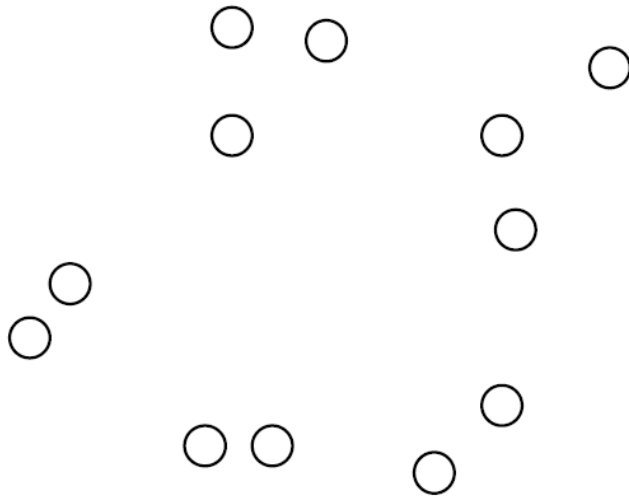
Clustering aglomerativ

- tehnica de grupare ierarhică mai populară
- algoritmul de bază este simplu
 - 1. Se calculează matricea de proximitate
 - 2. Fiecare punct de date va reprezenta un cluster
 - 3. Se repetă până când rămâne un singur cluster:
 - 4. Se unesc cele mai apropiate două cluster
 - 5. Se actualizează matricei de proximitate
- operaţiunea cheie este calcularea proximităţii a două cluster



Situaţia de pornire

- Se începe cu clustere de puncte individuale şi o matrice de proximitate



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

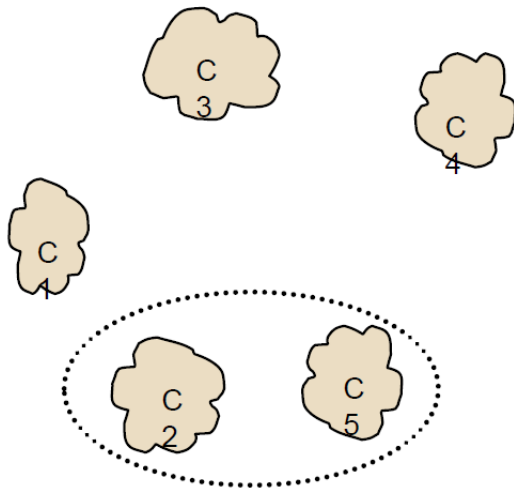
Proximity Matrix





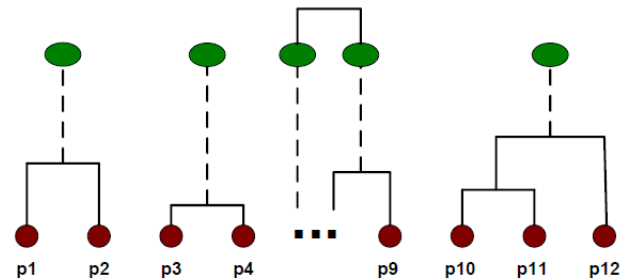
Situaţie intermediară

- după câteva etape de fuzionare, avem câteva clustere
- dorim să fuzionăm cele mai apropiate două clustere (C2 şi C5) şi să actualizăm matricea de proximitate



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

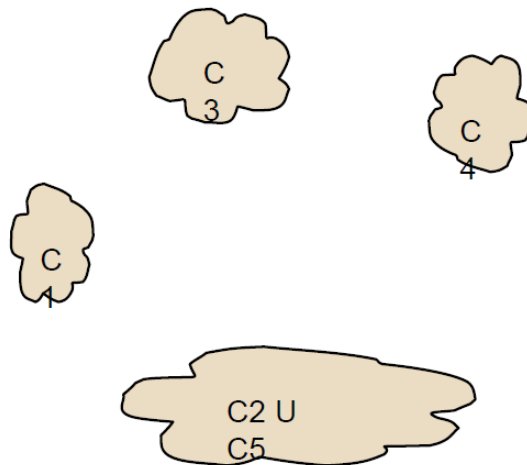
Proximity Matrix



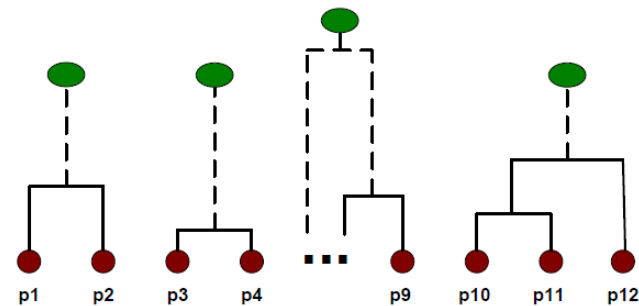


După fuzionare

- întrebarea este: "Cum actualizăm matricea de proximitate?"
adică "Cum măsurăm proximitatea (distanţa, similaritatea)
dintre două clusterere?"

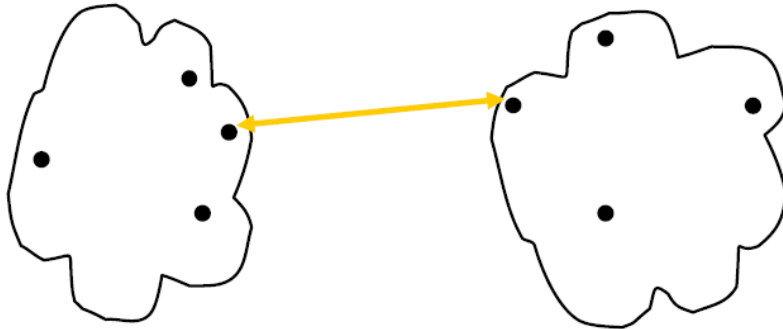


	C1	C2 U C5	C3	C4
C1		?		
C2 U C5	?	?	?	?
C3		?		
C4		?		

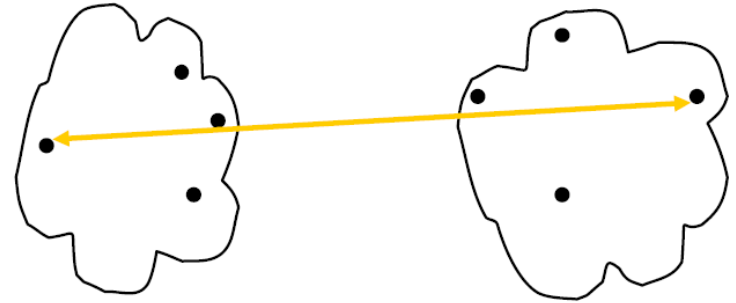




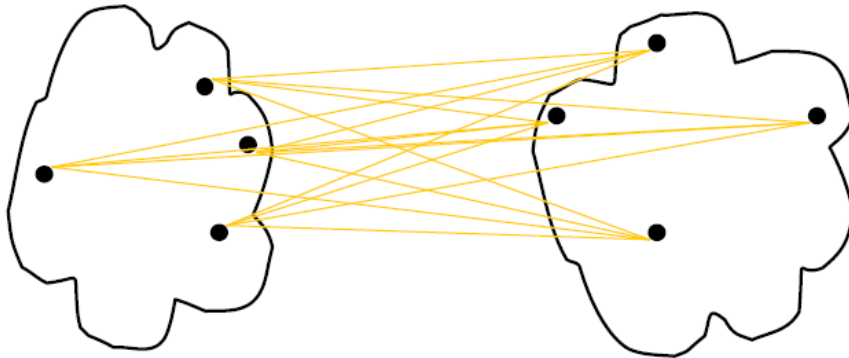
Similitudinea clusterelor



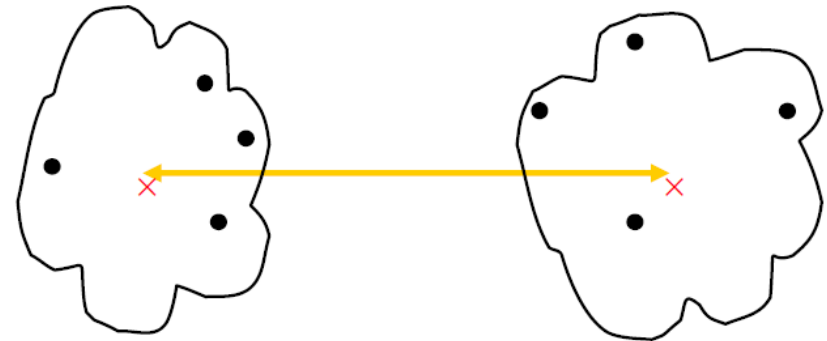
MIN (Legătură simplă)



MAX (Legătură completă)



Media grupului

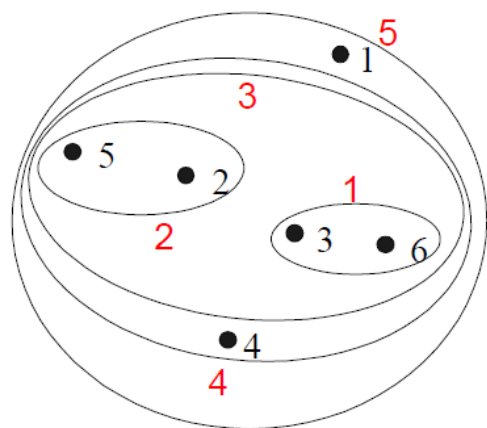


Distanţa dintre centroizi

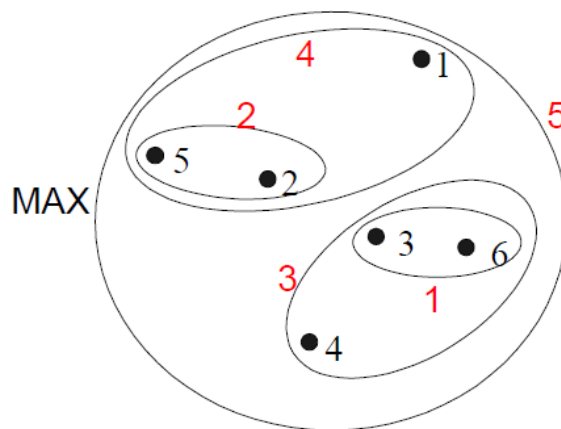
$$\text{proximity}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{proximity}(p_i, p_j)}{|\text{Cluster}_i| * |\text{Cluster}_j|}$$



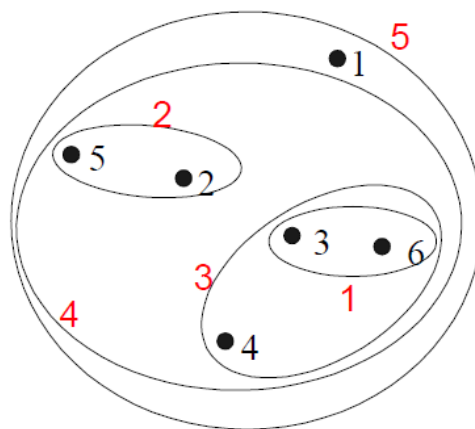
Comparaţie



MIN



MAX



Group Average





Clustering K-Means

- **K-means** este un algoritm partiţional de clusterizare
- setul de puncte de date (sau instanţe) D să fie

$$\{x_1, x_2, x_3, \dots, x_n\},$$

unde $x_i = (x_{i1}, x_{i2}, \dots, x_{ir})$ este un vector într-un spaţiu cu valori reale $X \subseteq R_r$, iar r este numărul de attribute (dimensiuni) din date

- algoritmul k-means împarte datele date în k cluster
 - fiecare cluster are un centru de cluster, numit ***centroid***
 - parametrul ***k*** este specificat de utilizator



Clustering K-Means

- Ideea de bază - inițializarea aleatorie a celor k centroide și iterarea prin pașii 2 și 3, până când centroidul nu se modifică:
 - 1. Se inițializează aleatoriu centroidele de cluster, c_1, \dots, c_k
 - 2. Având în vedere centroidele clusterelor, se determină punctele din fiecare grup:
 - pentru fiecare punct p , se găsește cel mai apropiat centroid c_i . Se introduce punctul p în clusterul i
 - 3. Considerând punctele plasate la pasul 2:
 - Se setează c_i ca fiind media punctelor din clusterul i
 - 4. Dacă s-a schimbat centroidul, se va repeta pasul 2



K-Means Clustering - Detalii

- centroizii inițiali sunt adesea aleși aleatoriu
 - clusterelor produse variază de la o execuție la alta
- centroidul este media punctelor din cluster
$$m_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$
- Distanța dintre puncte - măsurată prin distanța euclidiană
- Algoritmul K-Means va converge pentru măsurile de similaritate comune menționate mai sus
- cea mai mare parte a convergenței are loc în primele câteva iterații



Criteriu de oprire

■ Criteriul de convergenţă

- nu se mai produce nicio realocare a punctelor de date la clustere diferite
- nicio modificare (sau modificare minimă) a centroizilor
- scăderea minimă a sumei erorilor pătratică (*Sum of Squared Error* - *SSE*) - folosit pentru evaluarea K-Means:

$$SSE = \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} dist(\mathbf{x}, \mathbf{m}_j)^2$$

\mathbf{x} - este un punct de date în cluster

\mathbf{m}_j este centroidul clusterului



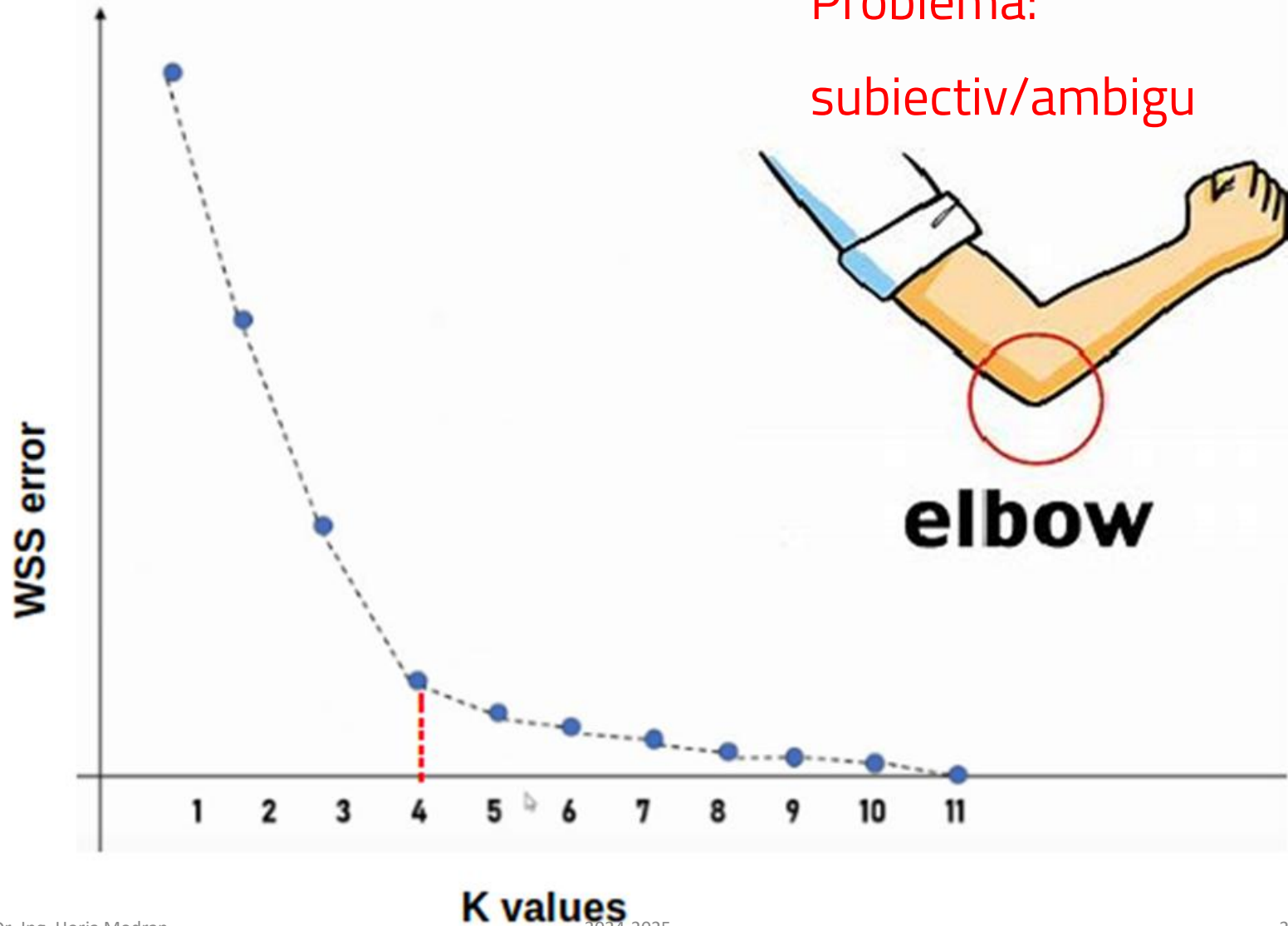
Metoda Elbow

- metoda "cotului" (*Elbow method*) este utilizată pentru a determina numărul de clustere dintr-un set de date
- se afişează variaţia explicată în funcţie de numărul de clustere
- alegerea "cotului" curbei ca număr optim de clustere
- **Inerţia** - suma distanţelor pătrate ale fiecărui punct de date faţă de cel mai apropiat centroid de cluster (eroarea pătrată totală a clusterizării)
- o valoare mai mică a inerţiei sugerează o grupare mai bună.



Metoda Elbow

Problemă:
subiectiv/ambigu





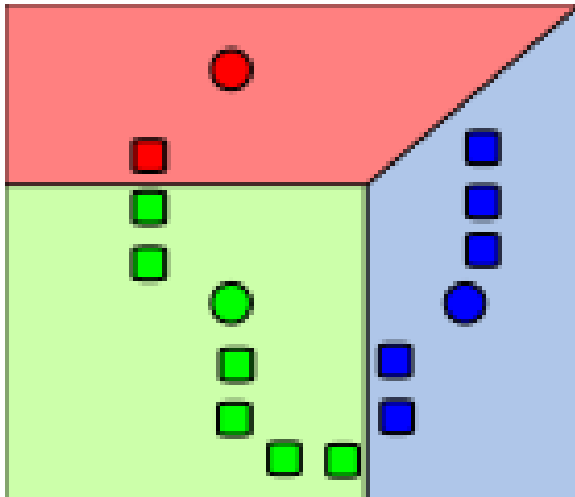
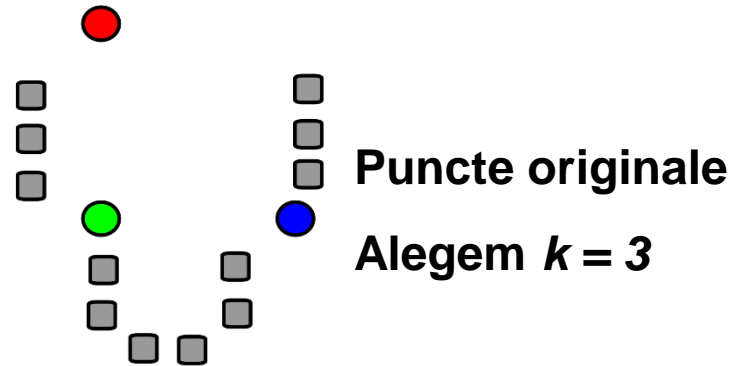
Coeficientul Silhouette

- măsoară cât de bine se încadrează un punct de date în clusterul atribuit în comparație cu alte cluster
- pentru un punct de date i , coeficientul **Silhouette** se calculează astfel:
$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

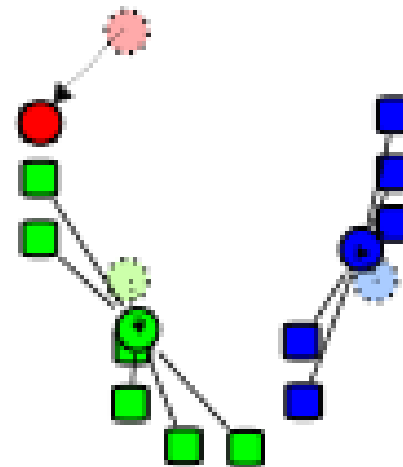
$a(i)$: distanța medie a lui i față de toate celelalte puncte din **propriul cluster** (distanța intracluster)
 $b(i)$: distanța medie a lui i față de toate punctele din cel **mai apropiat alt cluster** (distanța intercluster)
- Coeficientul **Silhouette** general (pentru tot algoritmul K-Means) este media valorilor $s(i)$ pentru toate punctele de date
- Interpretare:
 - +1: Punctul de date este atribuit corect la clusterul său
 - 0: Punctul de date se află în apropierea limitei de decizie
 - -1: Punctul de date este probabil clasificat greșit



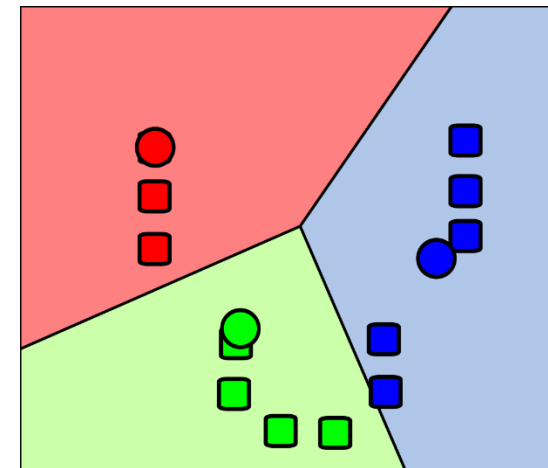
Exemplu



sunt create k cluster



Centroidul fiecăruia
dintre cele k cluster
devine noua medie

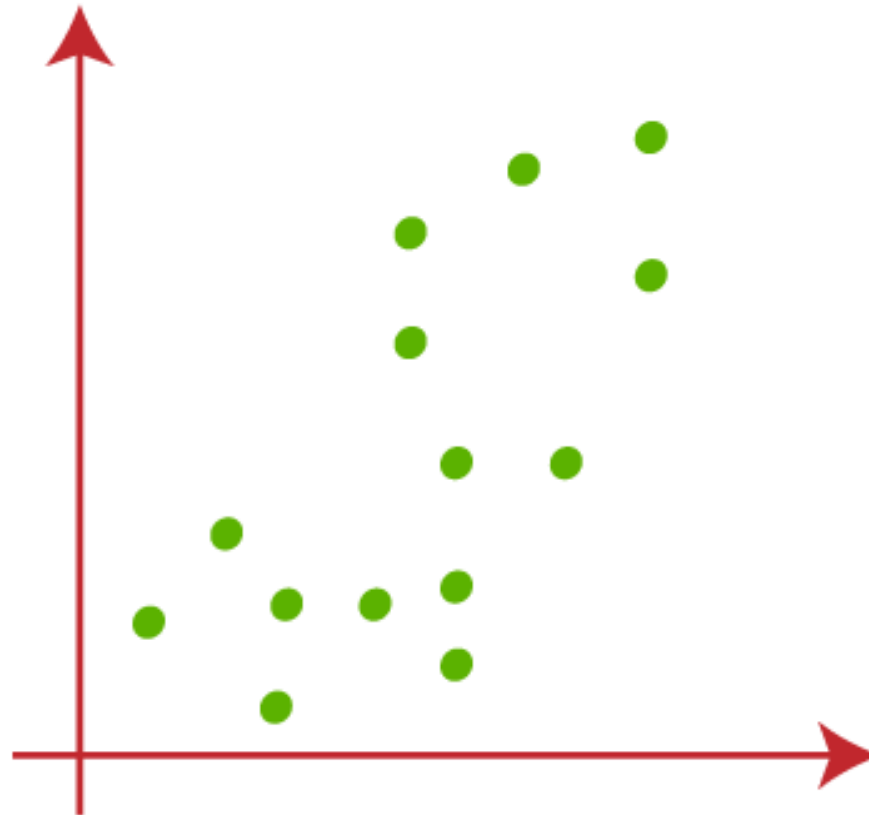


Etapele 2 şi 3 se
repetă până la
convergenţă



Exemplu

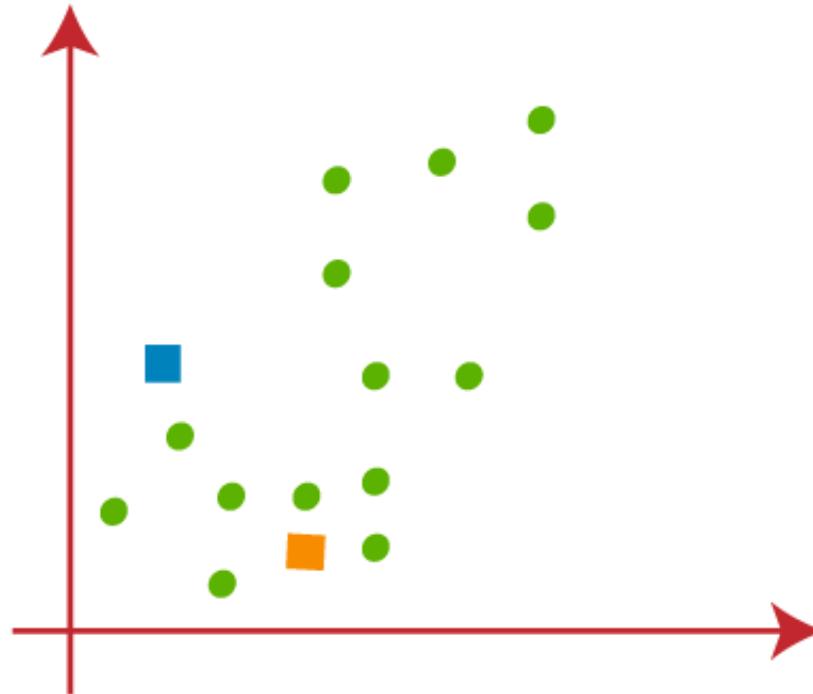
- Să presupunem că avem două variabile x şi y . Diagrama *scatter plot* pe axa x - y a acestor două variabile este ilustrată mai jos





Exemplu

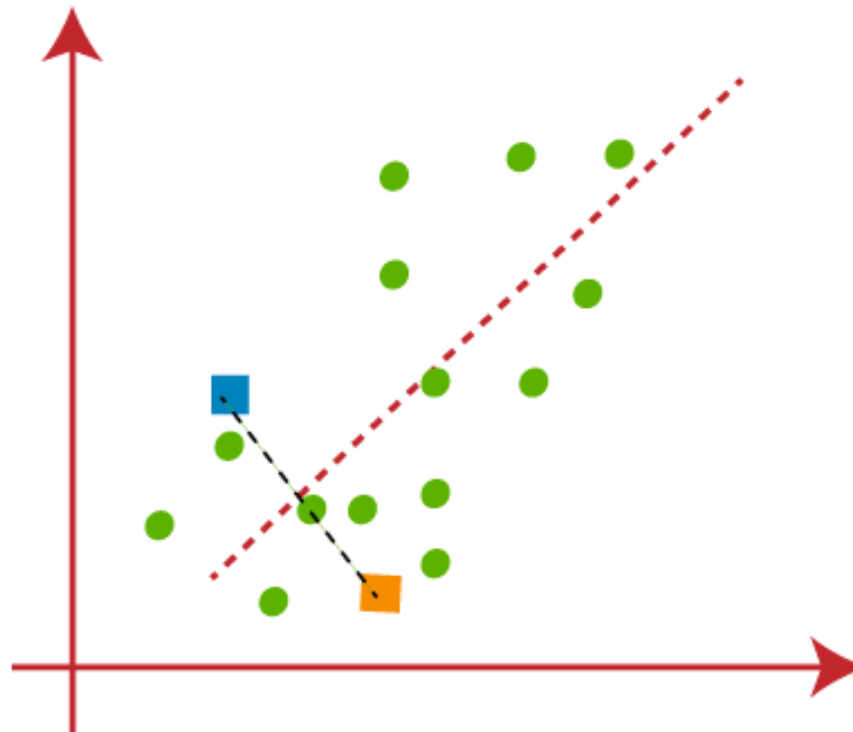
- Inițializăm $K=2$
- Trebuie să alegem aleatoriu k centroizi pentru a forma clusterelor (pot fi puncte din setul de date sau orice alte puncte)





Exemplu

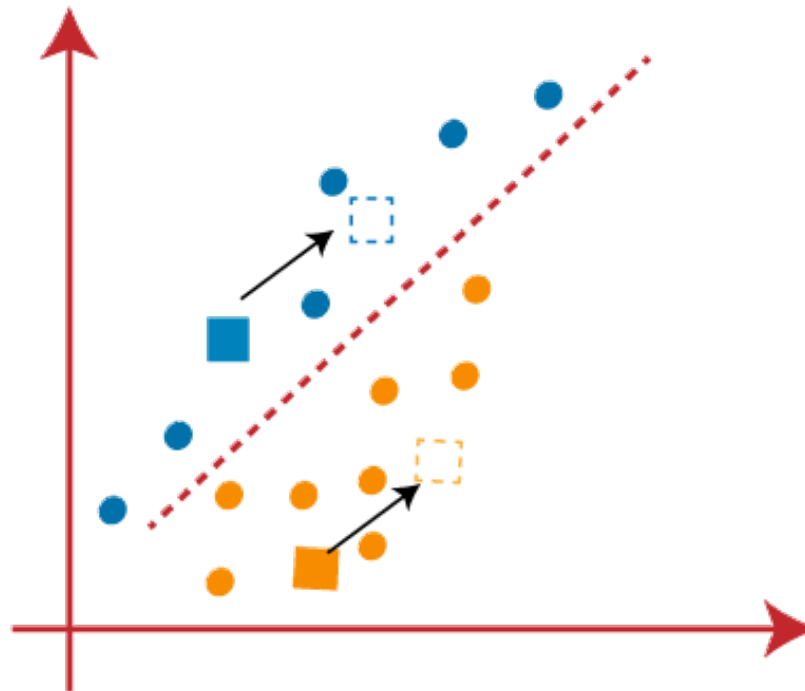
- Vom atribui fiecărui punct de date din grafic cel mai apropiat punct K (centroid)
- Trasăm o mediană între cei doi centroizi





Exemplu

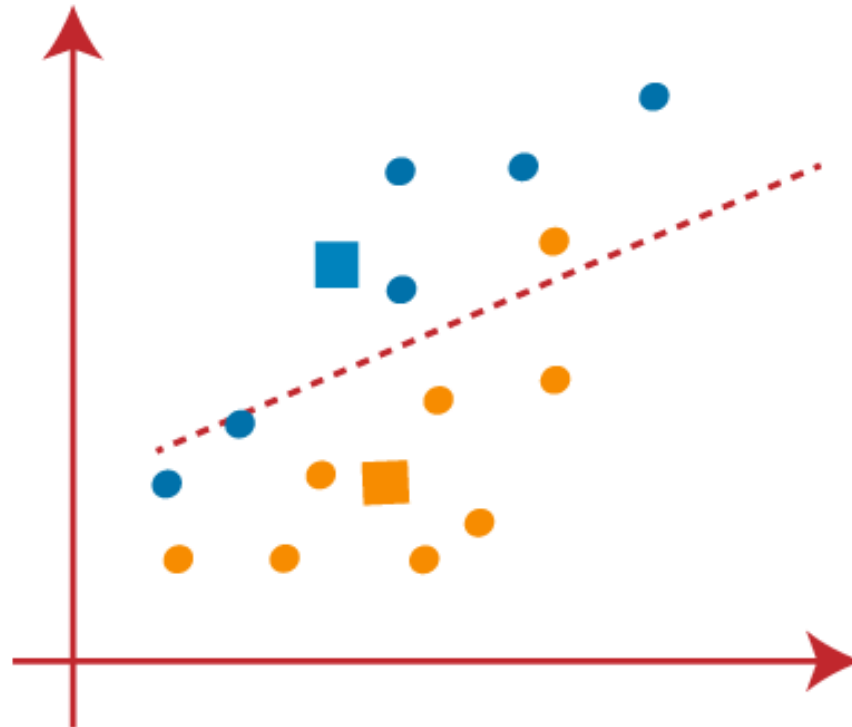
- Clusterele arată astfel
- Trebuie să găsim cel mai apropiat cluster pentru fiecare punct, deci vom determina noi centroizi





Exemplu

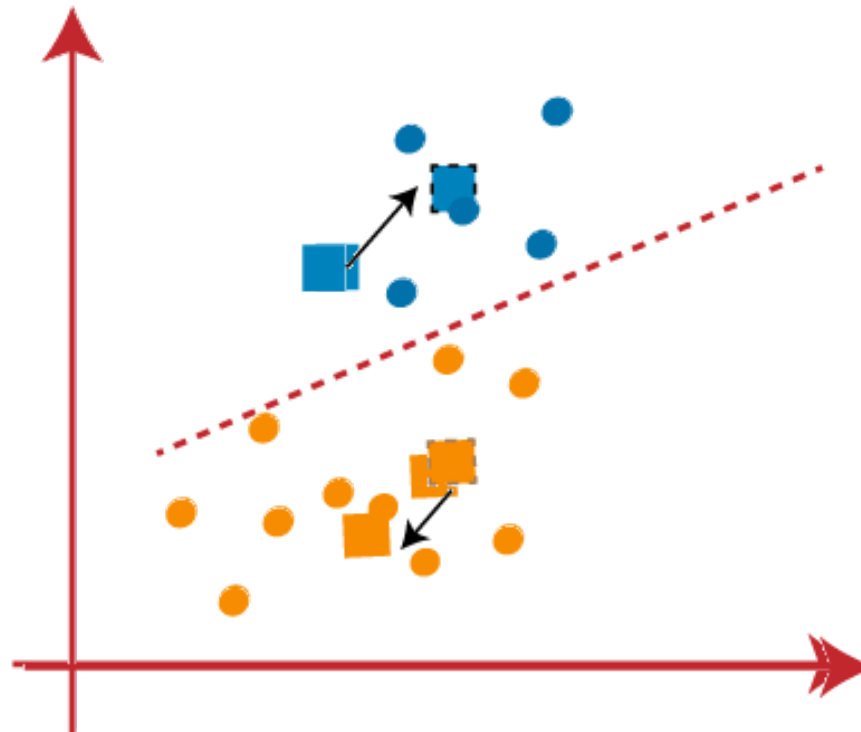
- Apoi, vom realoca fiecare punct de date noului centroid. Pentru aceasta, vom repeta acelaşi proces de găsire a unei linii mediane. Mediana este reprezentată punctat în graficul de mai jos:





Exemplu

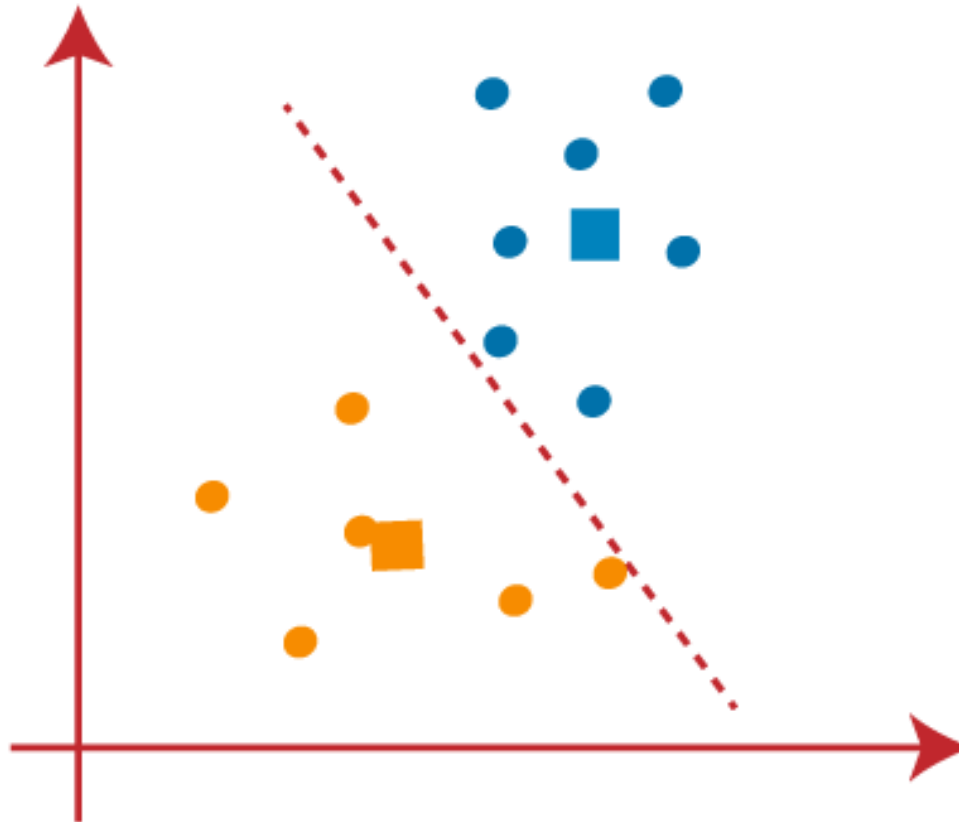
- Deoarece reassignarea a avut loc, vom trece din nou la pasul 3, care constă în găsirea unor noi centroizi





Exemplu

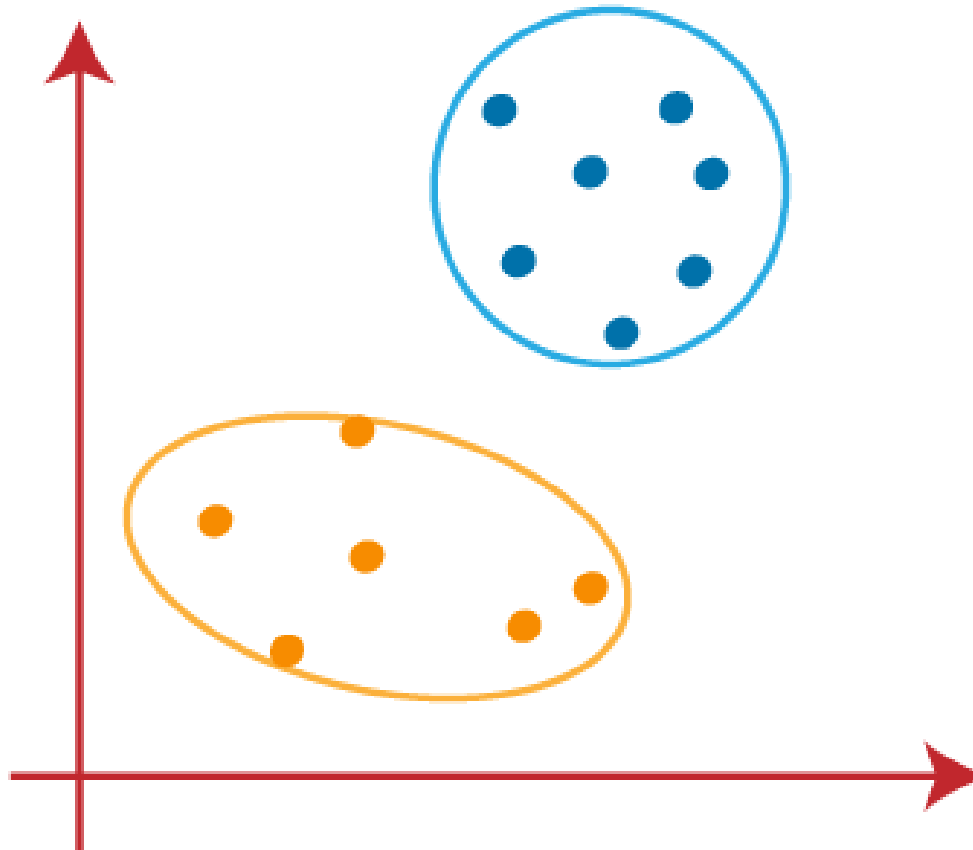
- Nu există puncte de date diferite pe fiecare parte a liniei, ceea ce înseamnă că algoritmul a ajuns la criteriul de oprire (convergență)





Exemplu

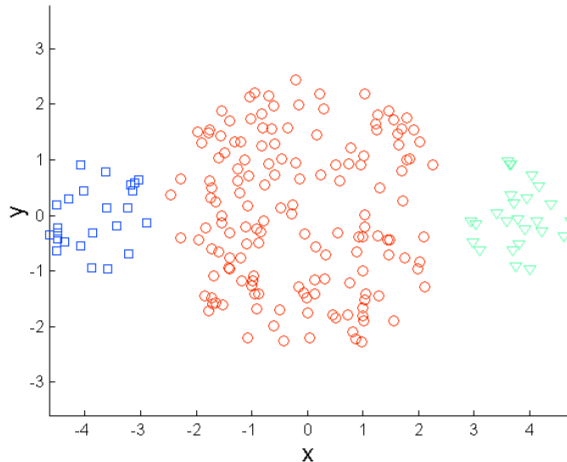
■ Clustere finale:



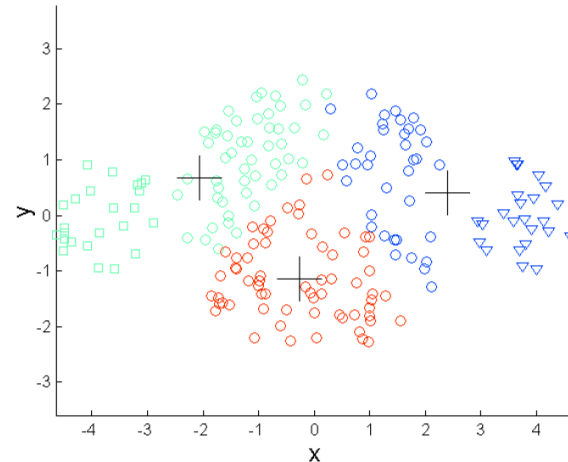


Limitări

- K-means are probleme atunci când clusterelor sunt de diferite
 - dimensiuni
 - densităţi
 - forme non-globulare
- K-means are probleme atunci când datele conţin *outliers*

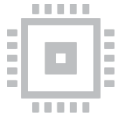


Original Points

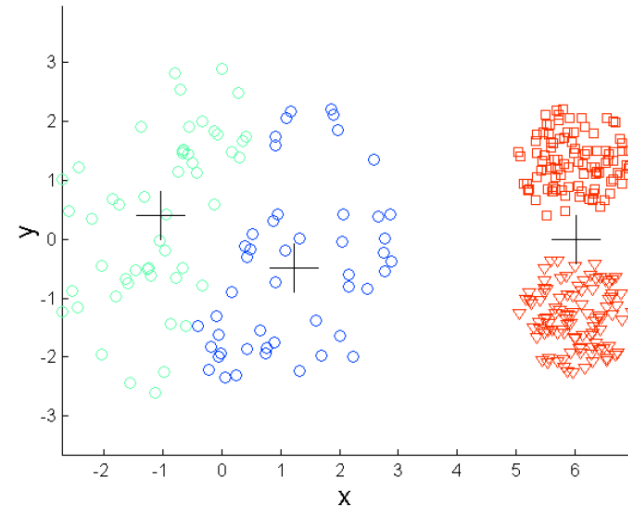
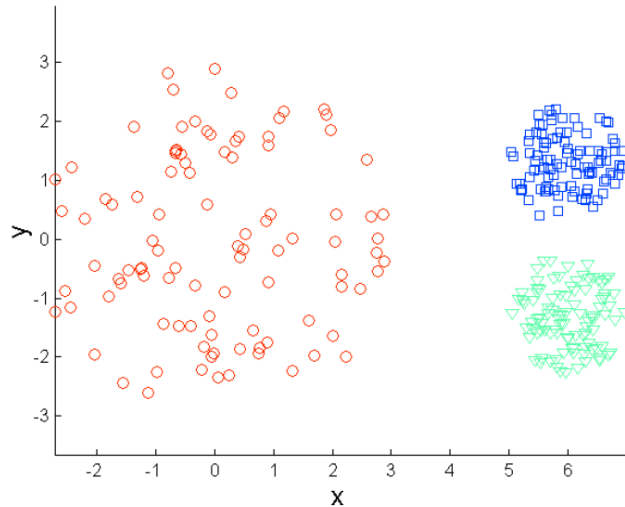


K-means (3 Clusters)

Dimensiuni
diferite

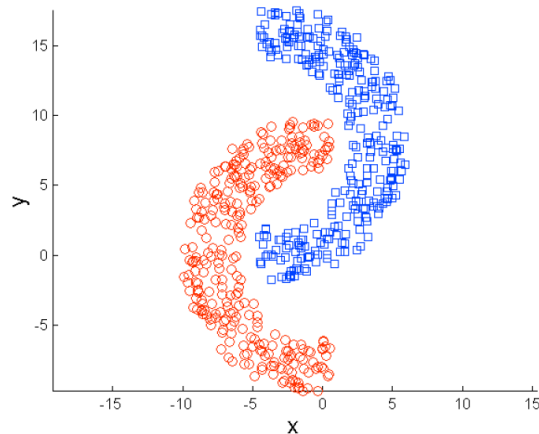


Limitări

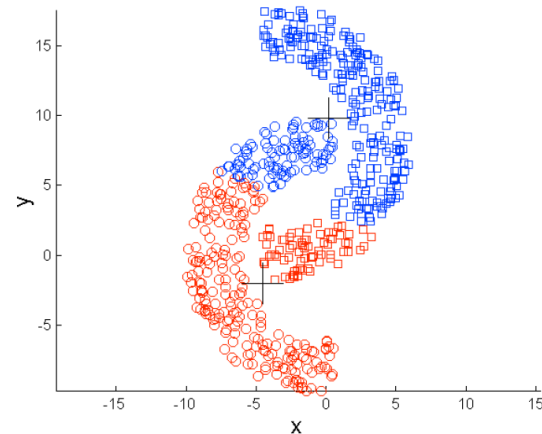


Densitate
diferită

Original Points



K-means (3 Clusters)



Forme non-
globulare



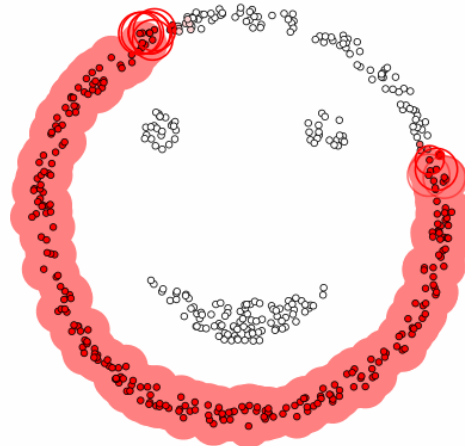
Avantaje şi dezavantaje

- Puncte forte:
 - simplu: uşor de înţeles şi de pus în aplicare
 - K-means este cel mai popular algoritm de clusterizare
 - se termină la un optim **local** dacă se utilizează SSE
- Slăbiciuni
 - nu funcţionează bine pentru diferite densităţi/forme
 - utilizatorul trebuie să specifice **k**
 - algoritmul este sensibil la valori de tip **outlier** (puncte de date care sunt foarte îndepărtate faţă de celelalte puncte)



DBSCAN

- Density Based Spatial Clustering of Applications with Noise(1996)
- un cluster este definit ca un set maxim de puncte conectate prin densitate
- descoperă cluster de formă arbitrară

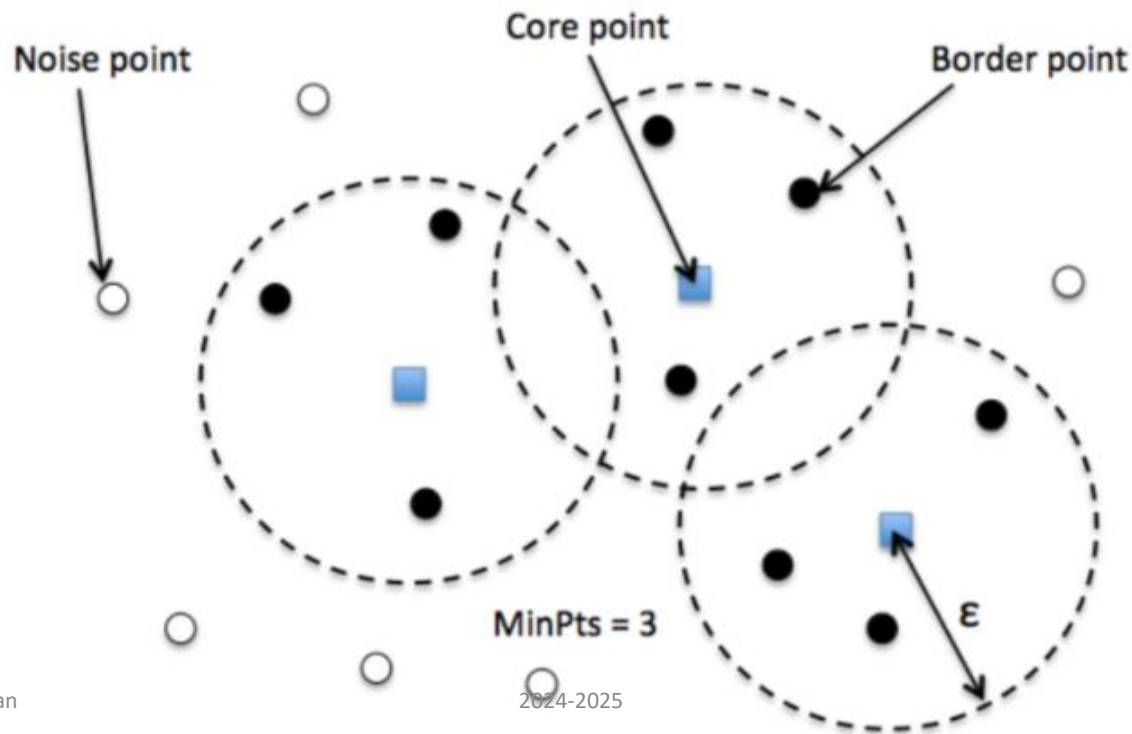




Tipuri de puncte

■ Tipuri de puncte

- **Core** – are cel puţin **MinPts** puncte la o distanţă ϵ faţă de el
- **Boder** – are cel puţin un punct Core la o distanţă ϵ
- **Outlier** – nu este nici core, nici border





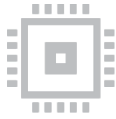
Algoritmul DBSCAN

- 1. Se selectează arbitrar un punct p
- 2. Se găsesc toate punctele accesibile din p în raport cu ε şi **MinPts**.
 ε - raza pentru vecinătatea punctului p
 $MinPts$ - numărul minim de puncte din vecinătatea dată $N(p)$
- 3. Dacă p este un punct **core**, se formează un cluster
- 4. Dacă p este un punct de tip **border**, niciun punct nu este accesibil din punct de vedere al densităţii din p şi DBSCAN vizitează următorul punct
- 5. Se continuă procesul până când toate punctele au fost procesate



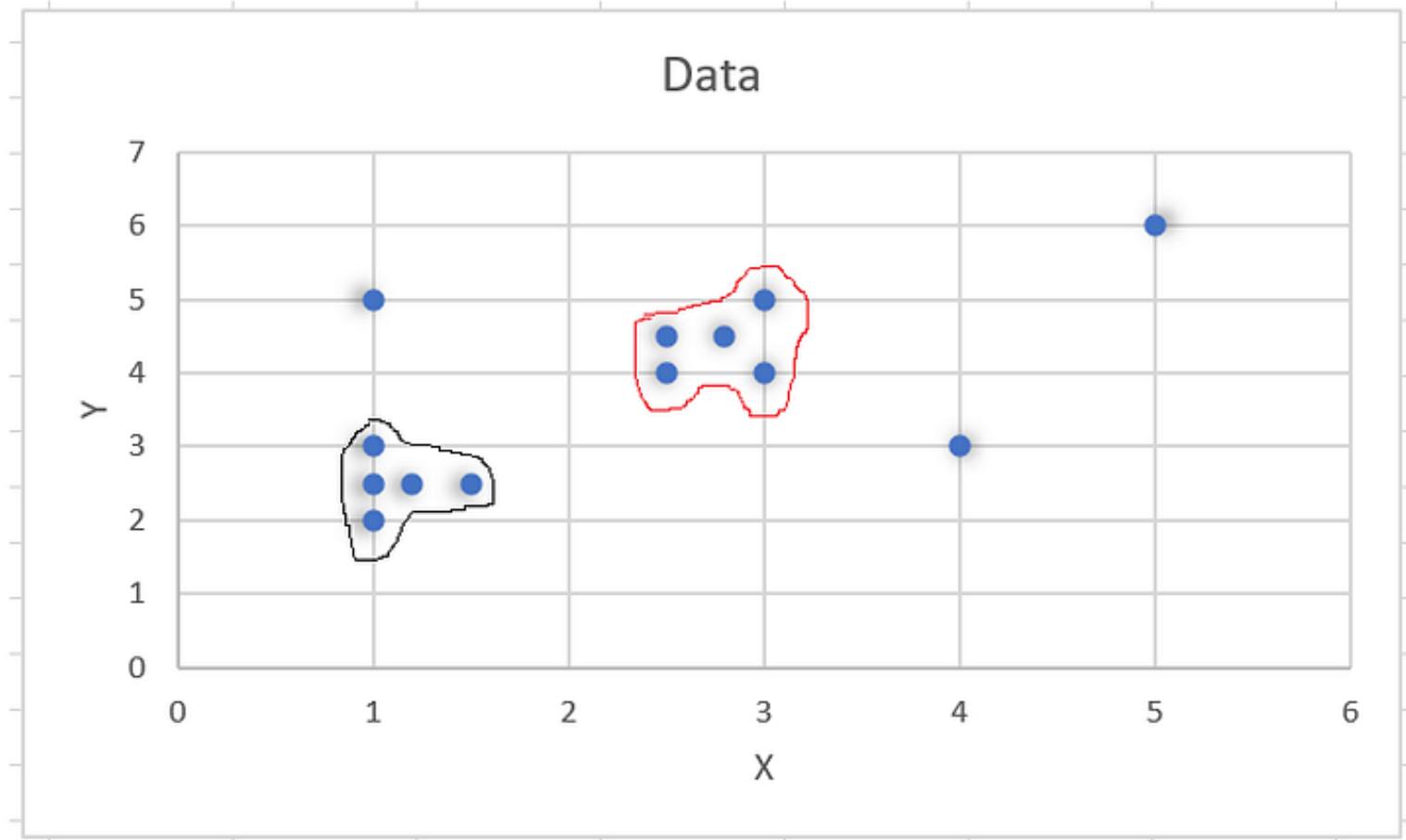
Estimarea parametrilor

- Pentru DBSCAN, sunt necesari parametrii ϵ şi **minPts**
- **minPts**: Ca regulă generală, $\text{minPts} \geq D + 1$ (de obicei $2 * D$)
 - valorile mai mari sunt de obicei mai bune pentru seturile de date cu zgomot
- dacă ϵ este ales mult prea mic, o mare parte din date nu vor fi incluse în niciun cluster, în timp ce pentru o valoare prea mare a lui ϵ clusterelor vor fuziona, iar majoritatea obiectelor vor fi în acelaşi cluster
 - în general, valorile mici ale lui ϵ sunt de preferat



DBSCAN - Exemplu

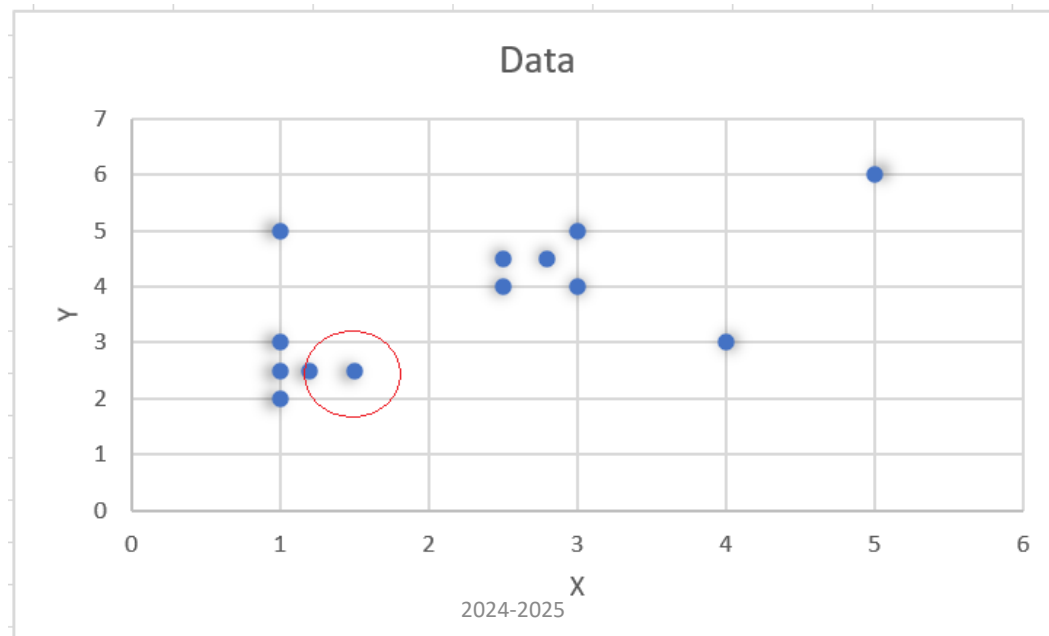
X	Y
1	2
3	4
2.5	4
1.5	2.5
3	5
2.8	4.5
2.5	4.5
1.2	2.5
1	3
1	5
1	2.5
5	6
4	3





DBSCAN - Exemplu

- Algoritmul DBSCAN are nevoie de doi parametri de intrare: raza din jurul fiecărui punct (ϵ) şi numărul minim de puncte de date care ar trebui să se afle în jurul aceluiaşi punct în raza respectivă (**MinPts**).
- De exemplu, considerăm punctul (1.5,2.5), dacă $\epsilon = 0.3$, atunci cercul din jurul punctului cu raza = 0.3, va conţine doar un alt punct în interiorul său (1.2,2.5), după cum se arată mai jos:





DBSCAN - Exemplu

- Alegem $\epsilon = 0.6$ şi **MinPts** = 4 şi considerăm primul punct de date din setul de date (1,2)

X	Y	Distance from (1,2)
1	2	0
3	4	2.8
2.5	4	2.5
1.5	2.5	0.7
3	5	3.6
2.8	4.5	3.08
2.5	4.5	2.9
1.2	2.5	0.53
1	3	1
1	5	3
1	2.5	0.5
5	6	5.6
4	3	3.1

Punctul (1, 2) are doar alte două puncte în vecinătatea sa (1, 2.5), (1.2, 2.5) - **este mai mic decât MinPts**



DBSCAN - Exemplu

- Repetăm procesul de mai sus pentru fiecare punct din setul de date şi determinăm vecinătatea fiecăruia. Calculele pot fi rezumate după cum urmează:

Point	Neighbourhood Points				
(1,2)	(1.2, 2.5)		(1, 2.5)		
(3, 4)	(2.5, 4)		(2.8, 4.5)		
(2.5, 4)	(3, 4)	(2.8, 4.5)	(2.5, 4.5)		
(1.5, 2.5)	(1.2, 2.5)		(1, 2.5)		
(3, 5)	(2.8, 4.5)				
(2.8, 4.5)	(3, 4)	(2.5, 4)	(3, 5)	(2.5, 4.5)	Cluster 1
(2.5, 4.5)	(2.5, 4)		(2.8, 4.5)		
(1.2, 2.5)	(1, 2)	(1.5, 2.5)	(1, 3)	(1, 2.5)	Cluster 2
(1, 3)	(1.2, 2.5)		(1, 2.5)		
(1, 5)					
(1, 2.5)	(1, 2)	(1.5, 2.5)	(1.2, 2.5)	(1, 3)	Cluster 2
(5, 6)					
(4, 3)					



DBSCAN - Exemplu

- Tabelul de mai jos arată clasificarea tuturor punctelor de date în puncte *core*, *border* şi *outlier*.

Point	Neighbourhood Points					
(1,2)	(1.2, 2.5)		(1, 2.5)		Border Point	
(3, 4)	(2.5, 4)		(2.8, 4.5)		Border Point	
(2.5, 4)	(3, 4)	(2.8, 4.5)	(2.5, 4.5)		Border Point	
(1.5, 2.5)	(1.2, 2.5)		(1, 2.5)		Border Point	
(3, 5)	(2.8, 4.5)				Border Point	
(2.8, 4.5)	(3, 4)	(2.5, 4)	(3, 5)	(2.5, 4.5)	Core Point	Cluster 1
(2.5, 4.5)	(2.5, 4)		(2.8, 4.5)		Border Point	
(1.2, 2.5)	(1, 2)	(1.5, 2.5)	(1, 3)	(1, 2.5)	Core Point	Cluster 2
(1, 3)	(1.2, 2.5)		(1, 2.5)		Border Point	
(1, 5)					Outlier	
(1, 2.5)	(1, 2)	(1.5, 2.5)	(1.2, 2.5)	(1, 3)	Core Point	Cluster 2
(5, 6)					Outlier	
(4, 3)					Outlier	



DBSCAN - Exemplu

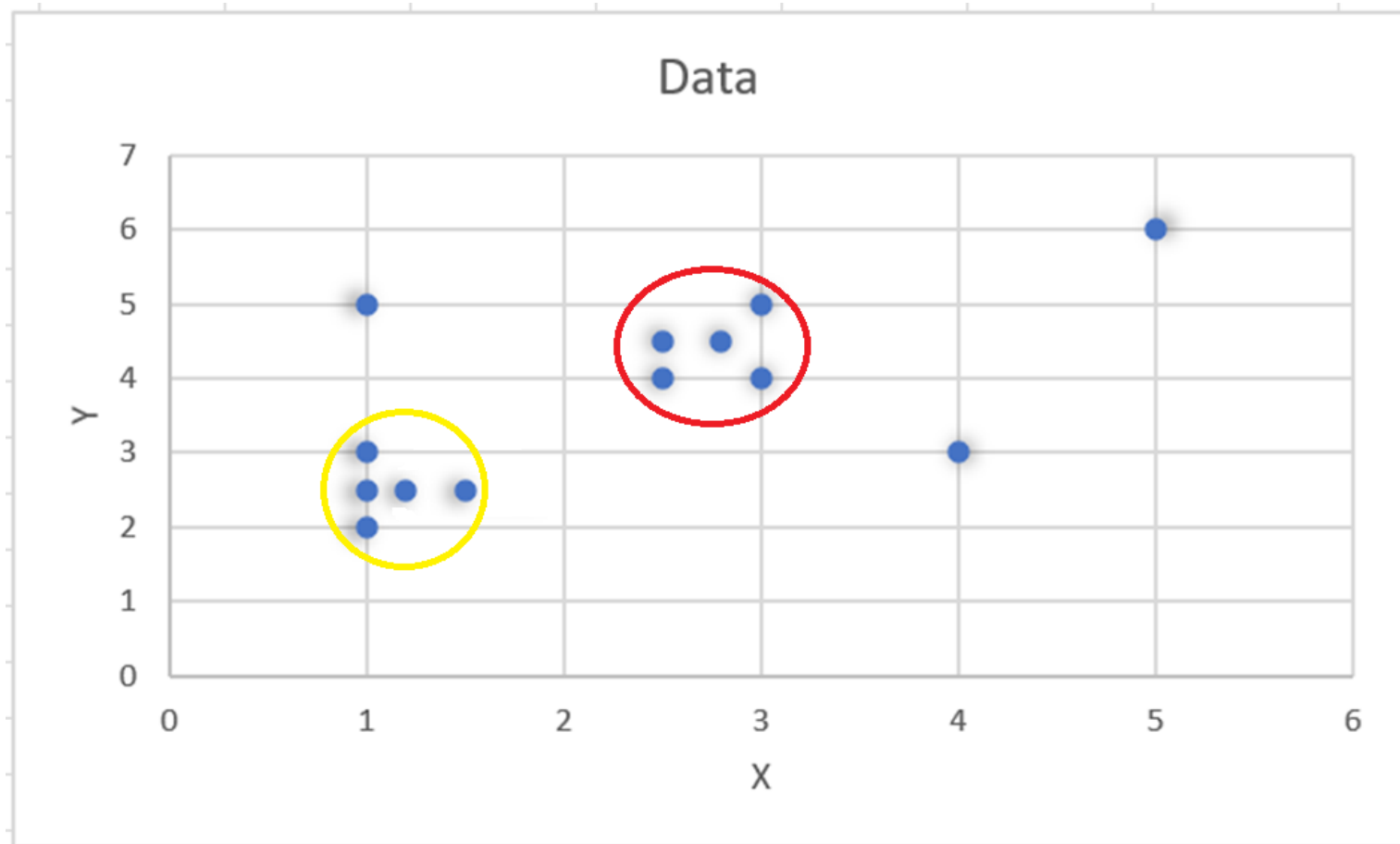
- Rezultatul algoritmului DBSCAN poate fi rezumat după cum urmează:

Cluster 1	Cluster 2	Outliers
(3,4)	(1, 2)	(1, 5)
(2.5, 4)	(1.5, 2.5)	(5, 6)
(3,5)	(1.2, 2.5)	(4, 3)
(2.8, 4.5)	(1, 3)	
(2.5, 4.5)	(1, 2.5)	



DBSCAN - Exemplu

■ Reprezentarea grafică a clusterelor rezultate:





ÎNTREBĂRI ?

