

Tour Version



Data Mining

# **Tour 5**

## **Clustering**

# Clustering Pattern

## clustering

What is a Cluster Pattern?

- Partitioning a set of Data into Subsets or Clusters

وهو ان أقسم البيانات لمجموعات كل مجموعة لها خصائص  
مشتركة يعني مثلاً الحيوانات لثدييات وطيور

- objects in a cluster are Similar yet  
Dissimilar to objects in other clusters

يعني ان ال objects اللى فى نفس ال Cluster متشابهين  
والعكس لو Cluster تاني

يعني مثلاً الطائر لازم كلهم بيبيضوا وعند هر ريش  
لوحبت كائن تاني مش كده يبقى مش طائر

Goal: discovery of previously unknown groups  
within the data  
- Implicit classes

الهدف هو إكتشاف أصناف جديدة من البيانات ما كناش واخدين  
بالناصها.

( أصناف مصنوعة - مفعية )

منظم یا مستخدمین Pre-Processing عملیات Remove noise

و کمالات outlier Analysis

## Requirements for Cluster Analysis

- Scalability  $\Rightarrow$  معظم الگوریتمها نمیتوانند Small Datasets را Handle کنند

فلو عندئذ Big Data را چطور تحلیل کنیم  
Sampling

- Handling Different attribute Types

معظم الگوریتمها فقط با Numerical Data کار میکنند

- Discovering clusters with arbitrary shape

معظم الگوریتمها فقط با شکل دایره‌ای کار میکنند

- Domain Knowledge & input Parameters

محتاج تفهیم اینته فایر! این من عملیات الگوریتم

و این هم الگوریتم الی Parameters الی مشخص الگوریتم



- Handling Noisy Data

حالياً أي Data فيها Noise يتأثر في شكل ال output يحتاج

- Incremental Clustering ( Sensitive )

موضوع الترتيب ديفرنس في عملية ال Clustering و محتاج تعديلها  
كنا مرة عشان تعرف حل النتيجة مستقرة ولا لا

- Handling High Dimensionality

معظم ال Algorithms تتعامل كمية قليلة من ال Attributes

- Constraint - Based Clustering

معظم ال Algorithms صاقد شرط واحد لها Constraint، ال حاجة  
بسطة زي ال K-mean ال بقدر بيس يتحكم من عدد ال Clusters

- Interpretability & Usability

هل ال Algorithms طالع نتائج مفهومة ونقدر نستخدمها .

## Comparing Cluster Analysis methods

الgorithms متکرم تقسما علی حسب

### 1. Partitioning Criteria

- Flat object که ~~cluster~~ بیتمی واحد
  - Hierarchical متکرم ال object بیتمی د
- Clusters different layers

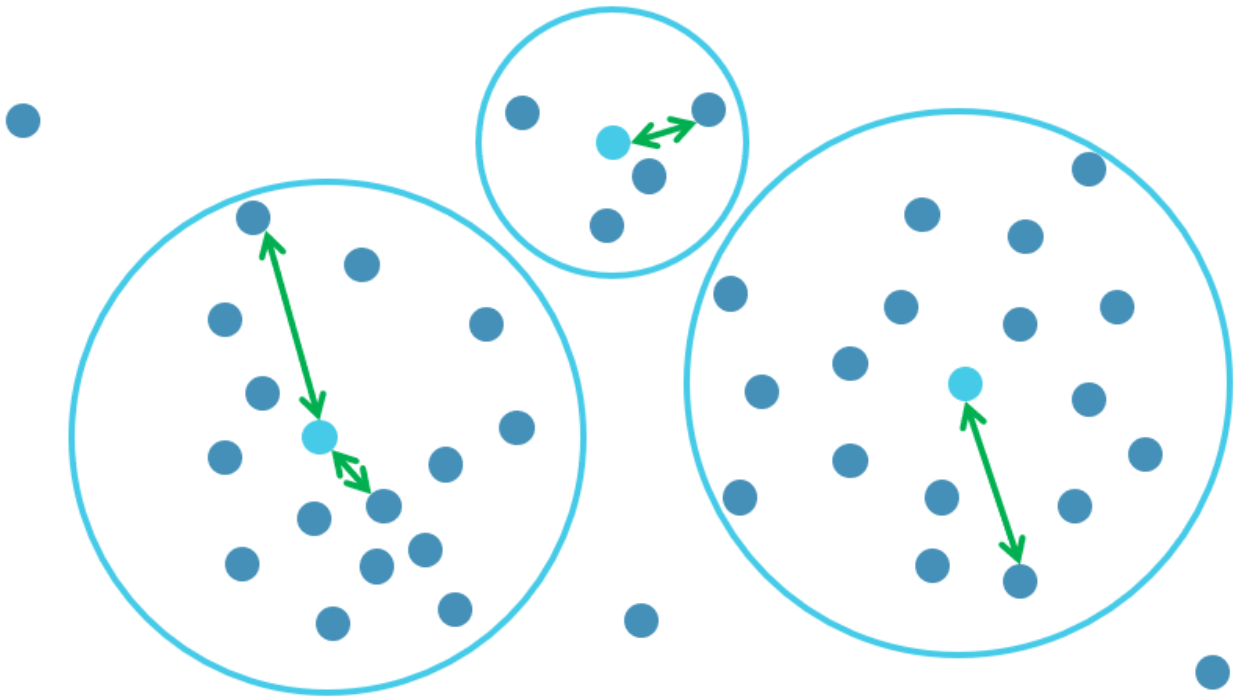
### 2. Separation of Clusters

صل که object بیتمی د Cluster واحد mutually exclusive  
از بیتمی زکتر (overlapping) Cluster

3. Similarity measure  
صل ال algorithms بیتمی ال distance و density Threshold

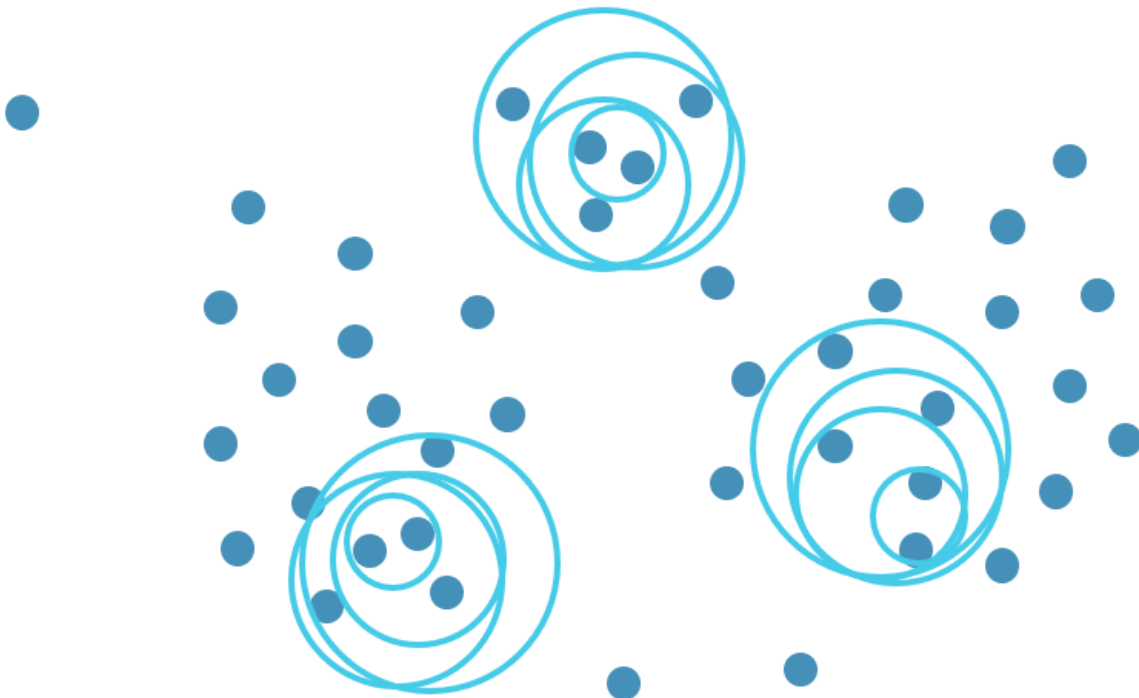
### Partitioning methods

- Find mutually exclusive clusters of spherical shape
- Distance-based
- May use mean or medoid to represent cluster center
- Effective for small- to medium-size data sets



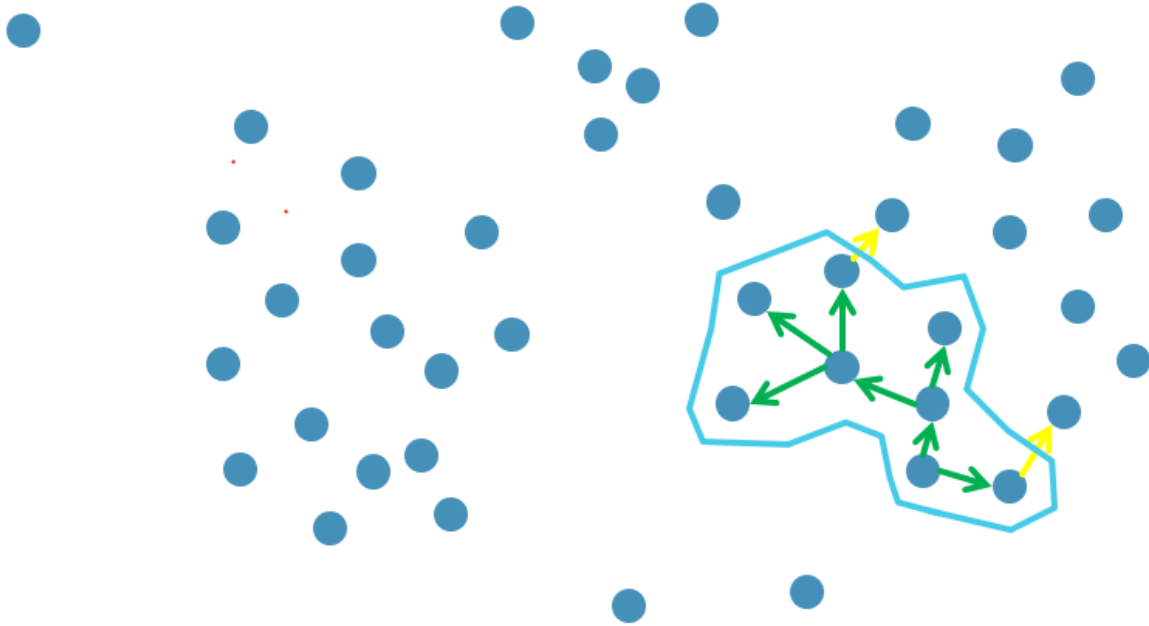
### Hierarchical methods

- Clustering is hierarchy involving multiple levels
- Cannot correct erroneous merges/splits
- May consider object "linkages"



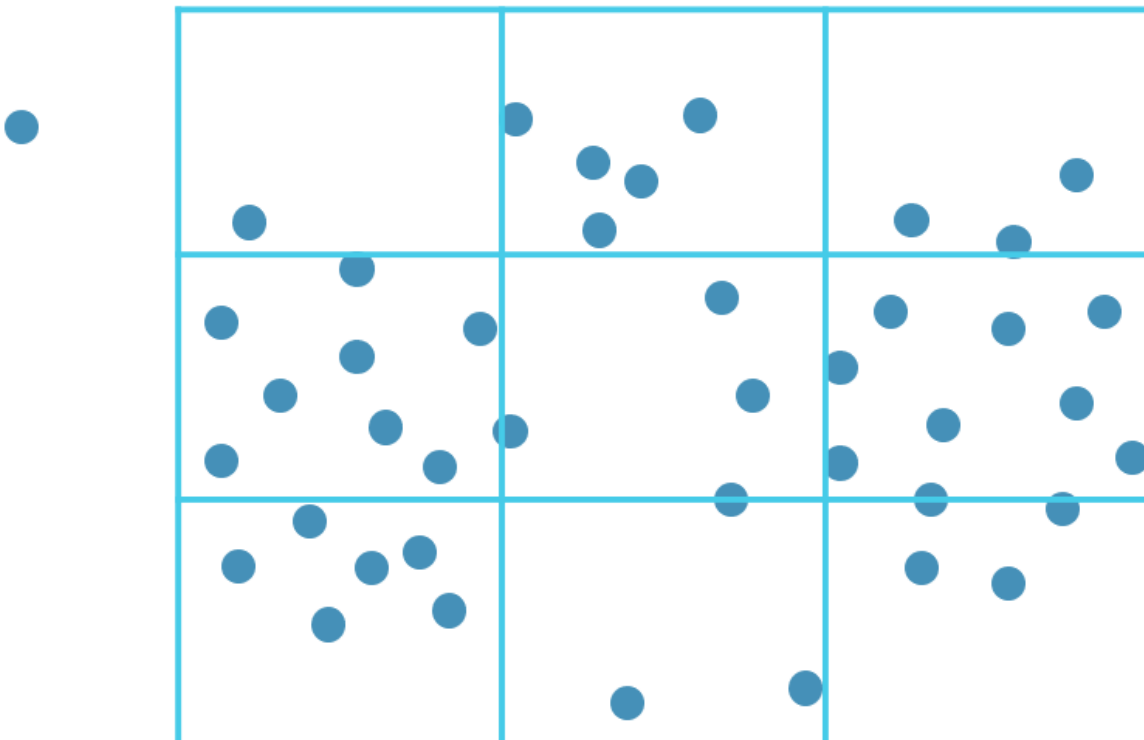
### Density-based methods

- Can find *arbitrarily shaped clusters*
- Clusters are *dense regions* separated by *low-density regions*
- Each point must have a *minimum number of points within its "neighborhood"*
- May *filter out outliers*



### Grid-based methods

- Use a multi-resolution *grid data structure*
- *Fast processing time*



## K-MEANS

يقسم البيانات لعدد  $k$  من ال Clusters  
كل Cluster التي بدورها هي نقطة المنتصف  
لل Cluster

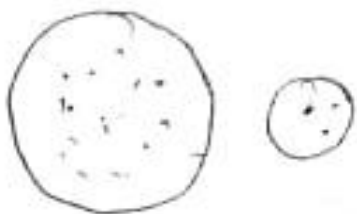
ال object الأقرب لل Centroid التي حددتهم يبعث تبع لأقرب  
Centroid من ال Cluster

Quality of Cluster  
حالة ان يكون الرقم ده أقل ما يمكن

$$E = \sum_{i=1}^k \sum_{x \in c_i} \text{dist}(x, c_i)^2$$

حاجات خدنا من اعتباراتك  
- اختيار عدد ال  $k$  كويس وايه من ال Centroids

• احسبه ال dissimilarity و ال mean كويس  
• لما تلاقي انت انجسام ال Clusters مختلفة جداً ارمش واحد  
Shape منتظم يبقى حصل Failure





	A1	A2
x1	2	10
x2	2	5
x3	8	4
x4	5	8
x5	7	5
x6	6	4
x7	1	2
x8	4	9

**Cluster the eight points in table using k-means.** Assume that  $k = 3$  and that initially the points are assigned to clusters as follows:  $C1 = \{x1, x2, x3\}$ ,  $C2 = \{x4, x5, x6\}$ ,  $C3 = \{x7, x8\}$ .

• Apply the k-means algorithm until convergence (i.e., until the clusters do not change), using the Manhattan distance.

(Hint: The Manhattan distance is:  $d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{in} - x_{jn}|$ .) Make sure you clearly identify the final clustering and show your steps.

**Algorithm: k-means.** The k-means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

**Input:**

- $k$ : the number of clusters,
- $D$ : a data set containing  $n$  objects.

**Output:** A set of  $k$  clusters.

**Method:**

- (1) arbitrarily choose  $k$  objects from  $D$  as the initial cluster centers;
- (2) repeat
  - (3) (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
  - (4) update the cluster means, that is, calculate the mean value of the objects for each cluster;
- (5) until no change;

K-MEANS

مبدأ خوارزمية K-MEANS عمل أول iteration

iter 1

$$C_1 = \{x_1, x_2, x_3\} \Rightarrow \bar{C}_1 = (4, 6\frac{1}{3}) \quad 3 = K \text{ عدد المجموعات}$$

$$C_2 = \{x_4, x_5, x_6\} \Rightarrow \bar{C}_2 = (6, 5\frac{2}{3})$$

$$C_3 = \{x_7, x_8\} \Rightarrow \bar{C}_3 = (2, 1\frac{1}{3})$$

iter 2	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	new mean
$C_1$	$5\frac{2}{3}$	$3\frac{1}{3}$	$6\frac{1}{3}$	$2\frac{2}{3}$	$4\frac{1}{3}$	$4\frac{1}{3}$	$7\frac{1}{3}$	$2\frac{2}{3}$	$(4\frac{1}{3}, 8\frac{1}{2})$
$C_2$	$8\frac{1}{3}$	$4\frac{2}{3}$	$3\frac{2}{3}$	$3\frac{1}{3}$	$1\frac{2}{3}$	$1\frac{2}{3}$	$1\frac{2}{3}$	$5\frac{1}{3}$	$(7, 4\frac{1}{3})$
$C_3$	5	1	7	5	5	5	5	5	$(1\frac{2}{3}, 5\frac{2}{3})$

$$\text{manhattan distance } (x_1, \bar{C}_1) = |4 - 2| + |10 - 6\frac{1}{3}| = 5\frac{2}{3}$$

$$" \quad " \quad (x_1, \bar{C}_2) = |6 - 2| + |10 - 5\frac{2}{3}| = 8\frac{1}{3}$$

$$" \quad " \quad (x_1, \bar{C}_3) = \dots = 5$$

$$C(x_1) = \min(5\frac{2}{3}, 8\frac{1}{3}, 5) = 5$$

وهكذا الفيت أما تدعى كل الجدول الذي فوقه

iter 2 clusters

$$C_1 = \{x_4, x_8\} \Rightarrow \bar{C}_1 = \frac{(5+9, 8+9)}{2} = (4\frac{1}{2}, 8\frac{1}{2})$$

$$C_2 = \{x_3, x_5, x_6\}$$

$$C_3 = \{x_1, x_2, x_7\}$$

- و هتفضل شغال كده في iterations لغيت اما تلاقي ال clusters شكلها ثبتت و ده الحل

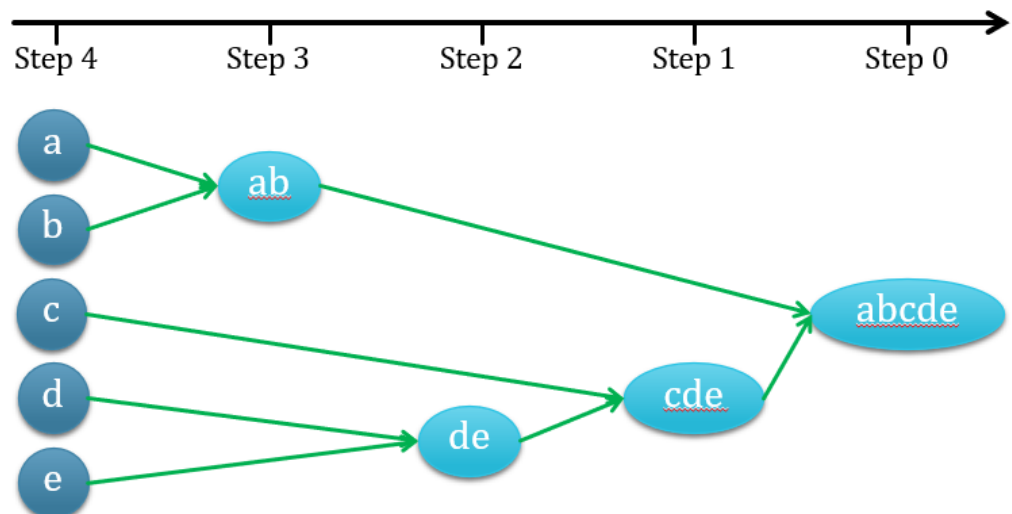
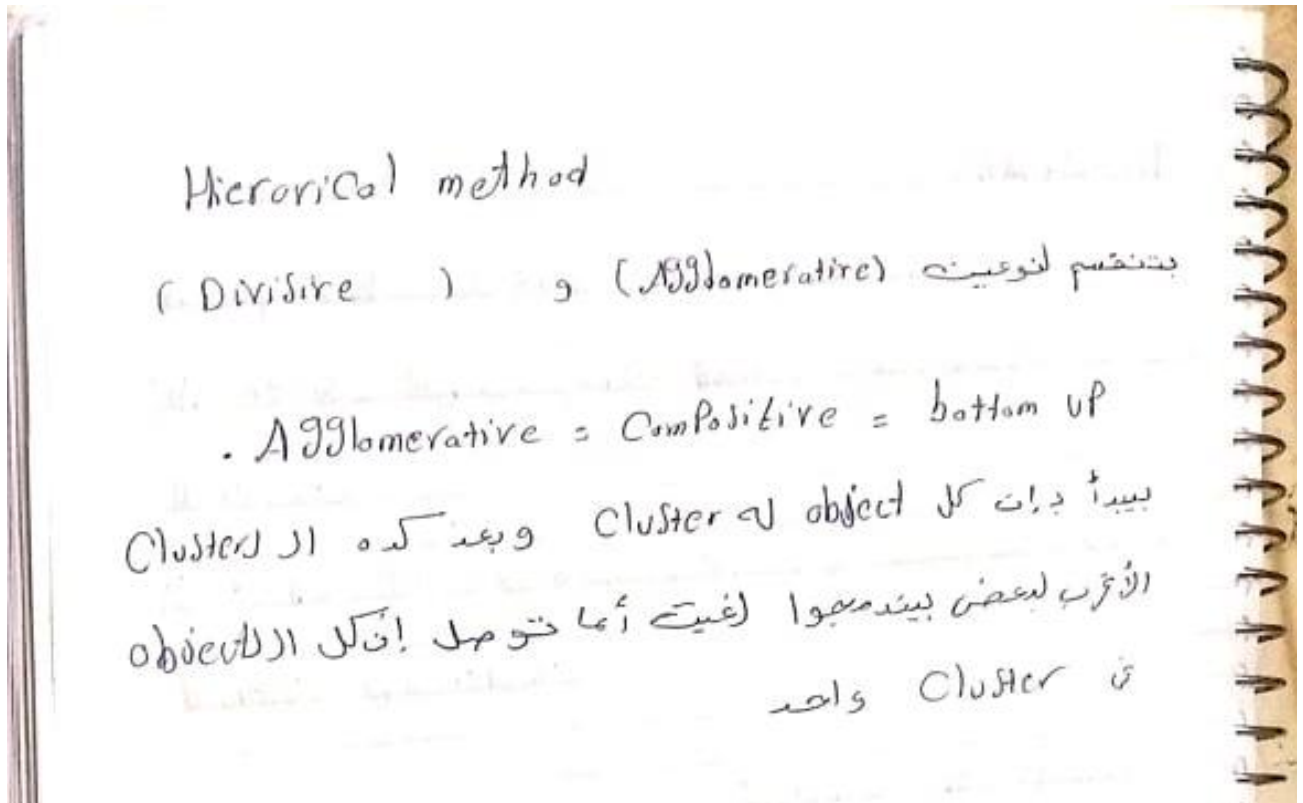
$C1 = \{x1, x4, x8\} = \{(2,10), (5,8), (4,9)\}$  Mean of  $C1 = (2\frac{2}{3}, 9)$

$C2 = \{x3, x5, x6\} = \{(8,4), (7,5), (6,4)\}$  Mean of  $C2 = (7, 4\frac{1}{3})$

$C3 = \{x2, x7\} = \{(2,5), (1,2)\}$  Mean of  $C3 = (1\frac{1}{2}, 3\frac{1}{2})$

## Hierarchical Methods {Agglomerative vs Devisive}

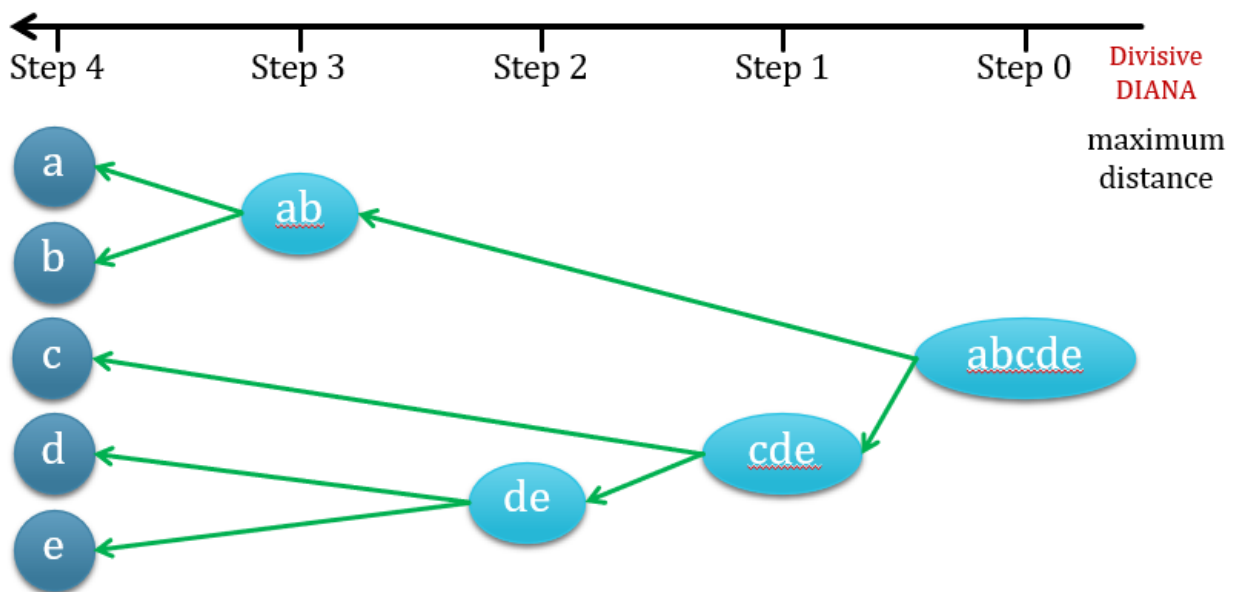
وينطلع فيها ال clusters على هيئة شجرة متتابعة و فيه منها نوعين من ال methods



Divisive = Top-down = Split decomposition

وہی مکی الی فاقہ پیدا ہے کہ اس (object) کی cluster

و بعد که بینقسموا لغیت اُمّا دکن لکن abcd Cluster



How to divide a cluster is a challenge! Heuristic approaches may be used

C



## Density Based Methods {DBSCAN}

DBSCAN

Density-Based Spatial Clustering of Applications with Noise

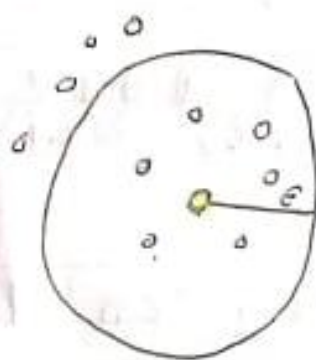
Find Core objects (with dense neighborhoods).

معطور [DB method] هو Core object کثافتی Partitioning

هو ال Centroid

بجاول! نه يلا Core obj يكون له Neighbors مصنفين كثافة

معيّن



Radius:  $\epsilon$

$\epsilon$ -neighborhood: all obj inside the circle

Min-PTJ

Threshold for Neighborhood

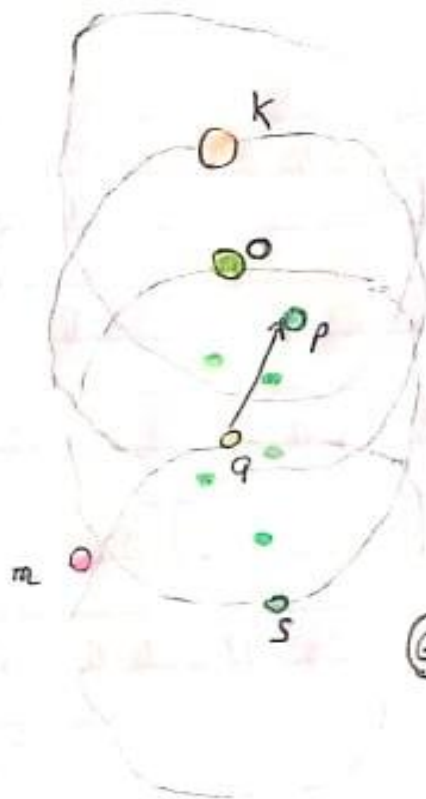
Density =  $\epsilon$  # of  $\epsilon$ -neigh

Core obj = ●

هو أي obj مصنف كثافة في ال region اللى تبعد أكثر من

أو بيساى Min-PTJ

Core object, obj whose  $\epsilon$ -neighborhood contains at least minPTJ objects



Given  $\epsilon = 4$  and  $\text{minPts} = 5$

③  $s$  متی Core لایه صفت Neighborhood

④  $q$  هو Core لایه صفت اکثر  $\epsilon = 4$  متی

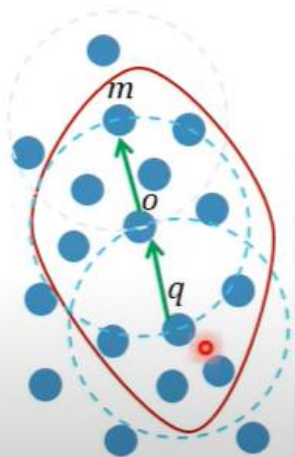
⑤  $p$  يعتبر Density Reachable لایه صفت  $q$  لایه صفت

⑥  $m$  متی Density Reachable لایه صفت  $q$  لایه صفت  $q$  لایه صفت

$q$  and  $k$  هم density connected لایه صفت عنصر مشترک لایه صفت

Core objects  $\Leftarrow q, k$  لایه صفت object بینشی لایه صفت Cluster 2 و

Given  $\epsilon = 4$  and  
MinPts = 5



an object  $p$  is directly **density-reachable** from another object  $q$  if and only if  $q$  is a core object and  $p$  is in the  $\epsilon$ -neighborhood of  $q$

objects  $q$  &  $m$  are **density-connected** if there is an object  $o$  such that  $q$  &  $m$  are both **density-reachable** from  $o$

---

**Algorithm:** DBSCAN: a density-based clustering algorithm.

**Input:**

- $D$ : a data set containing  $n$  objects,
- $\epsilon$ : the radius parameter, and
- $MinPts$ : the neighborhood density threshold.

**Output:** A set of density-based clusters.

**Method:**

- (1) mark all objects as **unvisited**;
- (2) **do**
- (3)     randomly select an **unvisited** object  $p$ ;
- (4)     mark  $p$  as **visited**;
- (5)     if the  $\epsilon$ -neighborhood of  $p$  has at least  $MinPts$  objects
- (6)         create a new cluster  $C$ , and add  $p$  to  $C$ ;
- (7)         let  $N$  be the set of objects in the  $\epsilon$ -neighborhood of  $p$ ;
- (8)         **for** each point  $p'$  in  $N$
- (9)             if  $p'$  is **unvisited**
- (10)                 mark  $p'$  as **visited**;
- (11)                 if the  $\epsilon$ -neighborhood of  $p'$  has at least  $MinPts$  points,  
                    add those points to  $N$ ;
- (12)             if  $p'$  is not yet a member of any cluster, add  $p'$  to  $C$ ;
- (13)         **end for**
- (14)         output  $C$ ;
- (15)     **else** mark  $p$  as **noise**;
- (16) **until** no object is **unvisited**;

منعمل Cluster Tendency عشاق نعره در Quality ال Clustering  
 تكتشف ان ال Data مالهاش شكل عشوائي

عن طريق Hopkins Statistics

← بنختار اول حاجة Sample عدد هم  $n$  من ال Points

وذلك نقطة من عشوائي ميت الا ترتب Neighbor لها ونحسب

ال distance ما بينها وما بين كل ال neighbor

$$\text{distance } X_i = \min_{v \in D} \{ \text{dist}(P_i, v) \}$$

نختار اقل سافة طاعت ذلك ال Samples

ونختار مربعة Sample عدد هم  $n$  ونعمل نفس ال عمل فوق  $\frac{1}{n}$

ونحسب النتيجة ونحسب

Hopkins Statistic

$H =$

$$\frac{\sum y_i}{\sum x_i + \sum y_i}$$

المفروض ان هى تساوى نص  $0.5$

$$H \approx 0.5$$

عشان ال method تكون اشتغلت مع



## Supervised Testing

Extrinsic method

مقارنة ما بين النتائج وأرض الواقع  
(Ground Truth)  
يطلبه ان يكون ما رغبت فعليا ان لكلا class  
بعدها ان لكلا object

والمقارنة ما بينهم ونشوف النسبة

$Q(C, C_g)$

محتاجين نعرف من خلاله

- Cluster Homogeneity

هل كل object له label class واحد بهي

على اساس ان Cluster بيمثل label class

- Cluster Completeness

بيشوف هل ان لكلا object له label class التي ان لكلا cluster يحتاجين

بهي من ان لكلا class

- Red Bag

ان لكلا object التي مالها ان لكلا cluster بنعينا فيها

- Small Cluster Preservation

من مثلا ان divisive hierarchical  
نرضي نحافظ على توازن الخصام

ان Clusters

BCubed

بالحسبة الـ Precision والـ Recall لكل object

Precision  $\Rightarrow$

كم عدد الـ objects التي في cluster واحد بيتقوا لنفس category

Recall  $\Rightarrow$

كم عدد الـ objects التي لهم نفس الـ label class لهم نفس cluster

**Extrinsic methods**  $\rightarrow$  compare clustering against **ground truth** (supervision)

○ Assign a score  $Q(C, C_g)$  to capture:

- **Cluster homogeneity**  $\rightarrow$  the purer the better – clusters represent separate class labels
- **Cluster completeness**  $\rightarrow$  an object with a class label belongs to the cluster representing that class label
- **Rag bag**  $\rightarrow$  objects that can't be merged into clusters belong to a *rag bag* – penalize a *misc. object* when put in a *pure cluster* more than in a *rag bag*
- **Small cluster preservation**  $\rightarrow$  splitting a small category is *more harmful* than splitting a large category

○ Ex. **BCubed** precision and recall of every object in dataset:

- Precision  $\rightarrow$  how many objects in the same cluster  $\in$  the same category as the object
- Recall  $\rightarrow$  how many objects of the same category are assigned to the same cluster

## Intuitive methods

وادي يتخبر أدب إلى ال Clusters التي طاعت مفصولة كويدي

مثال عليها Silhouette Coefficient

جنتفد رتفرق ما بين

① مقدار الفرق من المسافة ما بين <sup>متوسط</sup> cluster وكل ال clusters الى  
مشاء من فدرس ال cluster

كل ما كانت المتوسط أقل طبعاً نتيجة أصح

② مقدار أقل متوسط لفرق المسافة ما بين ال clusters وباقي

ال cluster دس مش ال cluster الى بيقتس له cluster

كل ما كانت كبير طبعاً هي يبقى أفضل

لو المقدار الثاني - الأول طالع موجب بيقتس cluster كويدي

لو العكس يبقى ال cluster عليها

•Compute **the silhouette coefficient** for object **x1**.

What is the meaning of the computed value?

$$C1 = \{x1, x4, x8\} = \{(2,10), (5,8), (4,9)\} \quad \text{Mean of } C1 = (2\frac{2}{3}, 9)$$

$$C2 = \{x3, x5, x6\} = \{(8,4), (7,5), (6,4)\} \quad \text{Mean of } C2 = (7, 4\frac{1}{3})$$

$$C3 = \{x2, x7\} = \{(2,5), (1,2)\} \quad \text{Mean of } C3 = (1\frac{1}{2}, 3\frac{1}{2})$$

$$a(o) = \frac{\sum_{o' \in c_i} \text{dis}(O, O')}{|c_i| - 1} = \frac{5+3}{2} = 4$$

$$b(o) = \min \left\{ \frac{\sum_{o' \in c_j} \text{dis}(O, O')}{|c_j|} \right\} = \min \left\{ \frac{12+10+10}{3}, \frac{5+9}{2} \right\} = 7$$

$$S(o) = \frac{b(o) - a(o)}{\max\{a(o), b(o)\}} = \frac{7-4}{7} \rightarrow +ve$$

○ This mean the cluster containing o is compact and o is far from other cluster