

Happy Friday!

- Project Proposal Due Tonight
- Peer Evaluation Out tonight
 - Due on Apr 20th

**THE CLUSTERS,
THEY SPEAKS TO US.**

Kmeans - Hard - K

GMM - Soft - K

Hierarchical Clustering - Hard - X

introduce hyperparameters
Single/Complete/
Centroid/Average
L2 Manhattan

DBSCAN - Hard - X


Density-Based Clustering

Nimisha Roy

Lecturer, College of Computing

Director, Online Undergraduate Initiatives

Outline

- Overview
- Basic Concepts 
- The DBSCAN Algorithm
- Analysis of DBSCAN

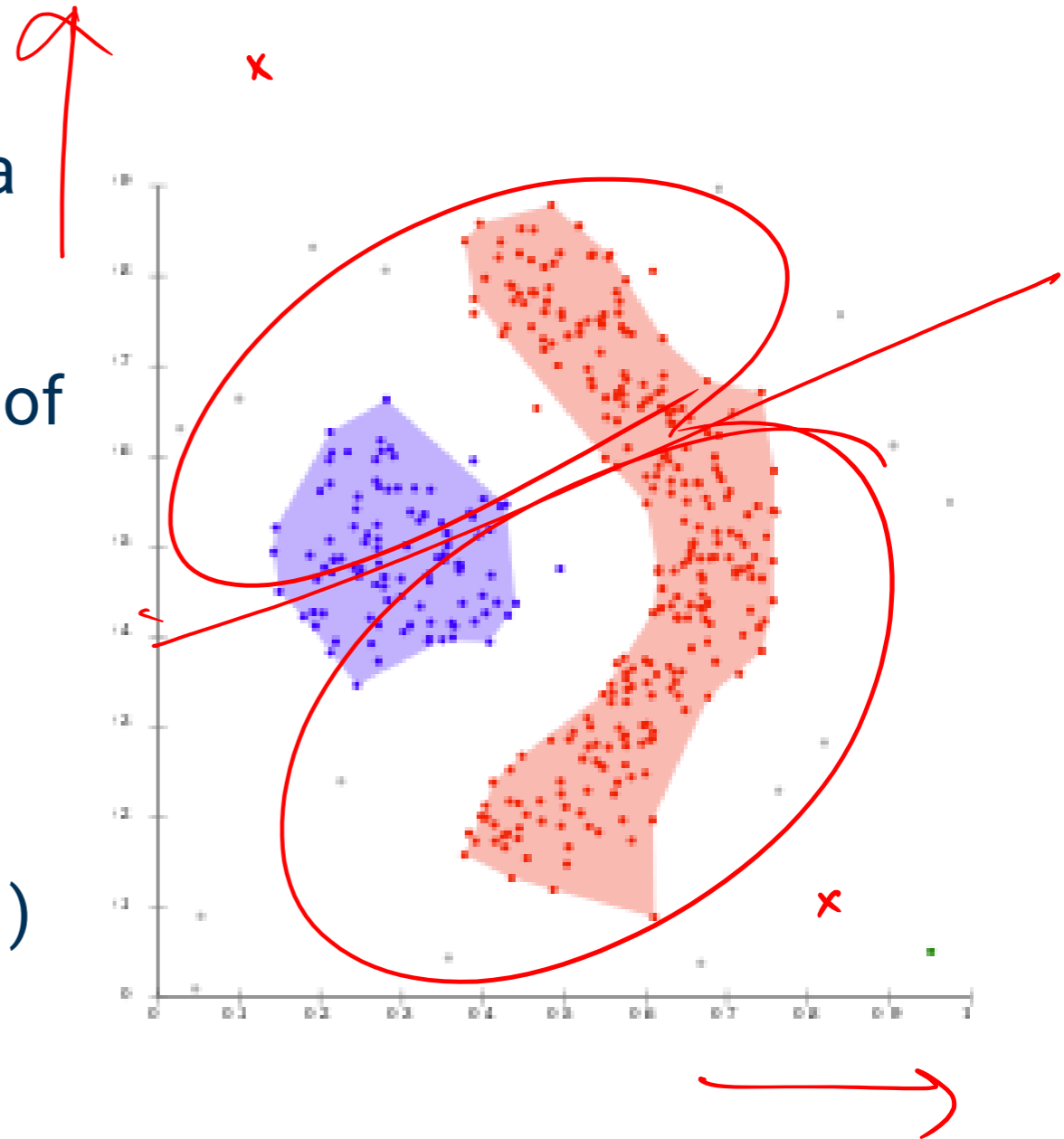
Density-Based Clustering

- Basic Idea

- Clusters are dense regions in the data space, separated by regions of lower density
- A cluster is defined as a maximal set of density-connected points
- Detect arbitrarily shaped clusters

- Method

- DBSCAN (Density-Based Spatial Clustering of Applications with Noise)



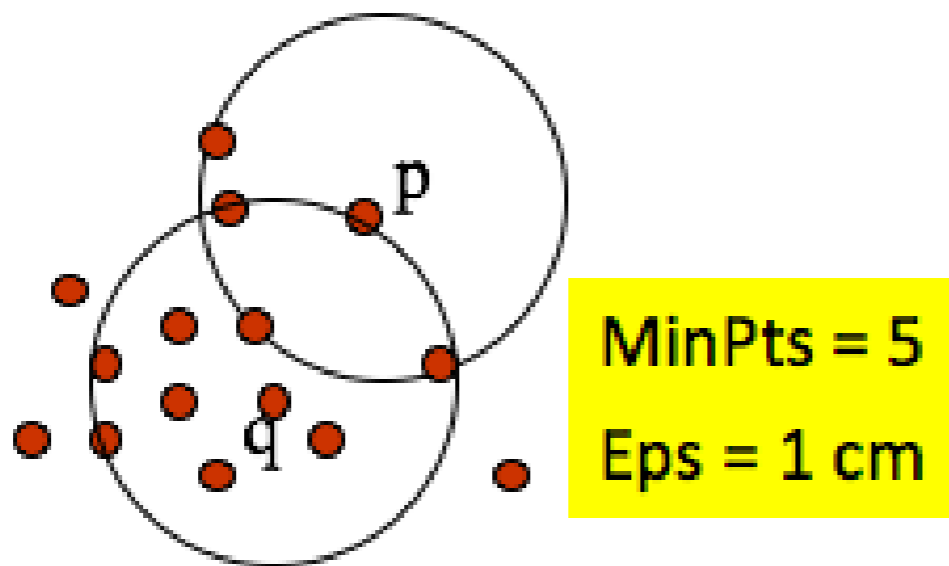
Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN



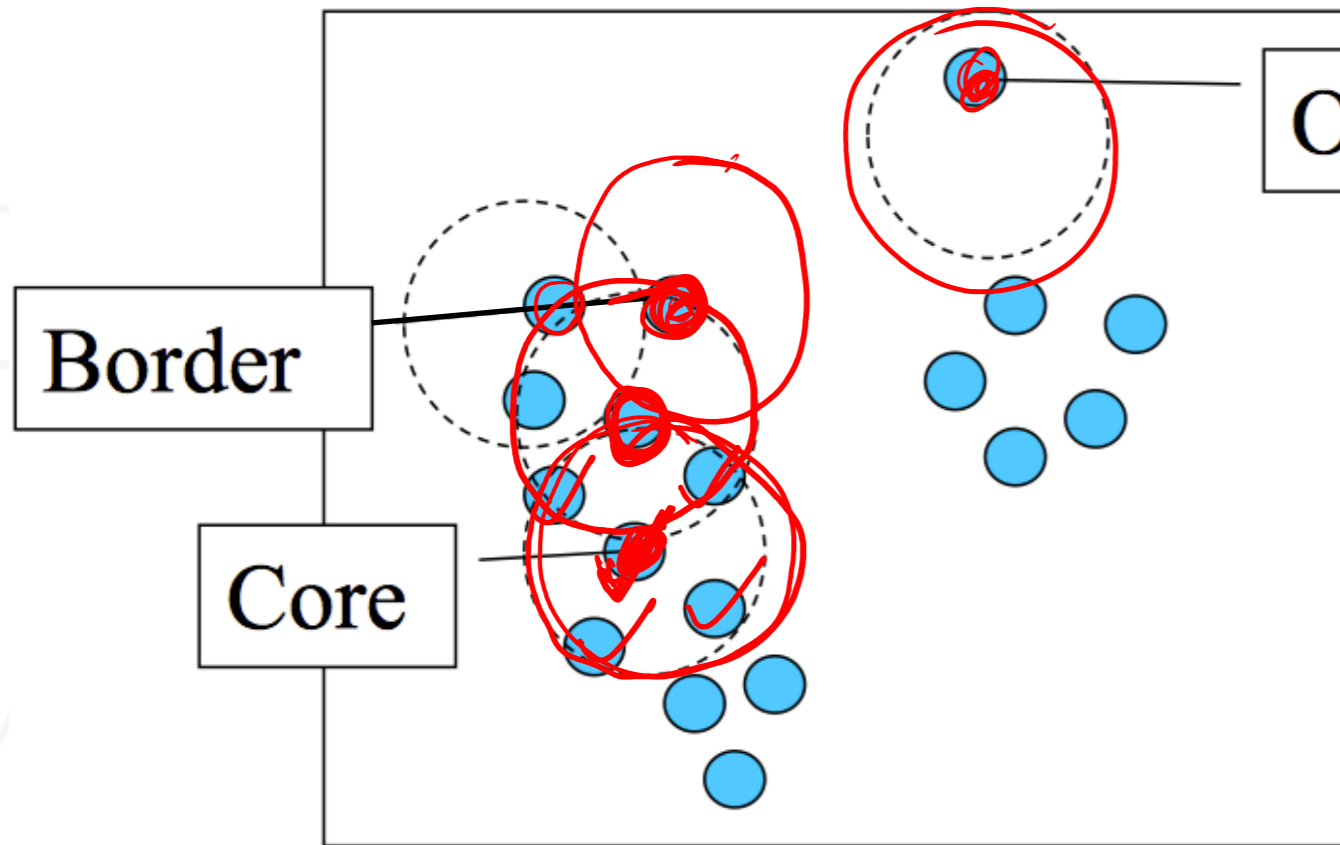
High Density v.s. Low Density

- Two parameters
 - **Eps (ϵ)**: Maximum radius of the neighborhood
 - **MinPts**: Minimum number of points in the Eps-neighborhood of a point
- High density: ϵ -Neighborhood of an object contains at least MinPts of objects



Density of p is low
Density of q is high

Core Points, Border Points, and Outliers



$\epsilon = 1 \text{ unit}$, $\text{MinPts} = 5$

Outlier

Given ϵ and *MinPts*, categorize the objects into three exclusive groups.

A point is a **core point** if it has more than a specified number of points (MinPts) within Eps—These are points that are at the interior of a cluster.

A **border point** has fewer than MinPts within Eps, but is in the neighborhood of a core point.

A **noise point** is any point that is not a core point nor a border point.

DBSCAN algorithm

Goal: Find clusters of arbitrary shapes based on density while marking noisy points as outliers.

Steps:

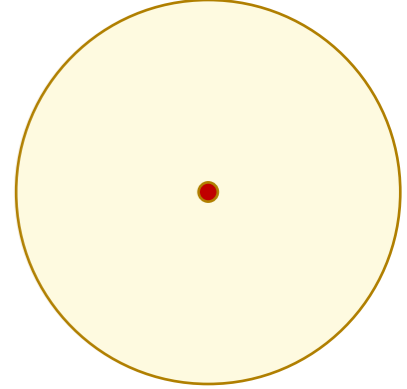
- **Pick any unvisited point.**
 - Mark it as *visited* so we don't process it twice.
 - If it has enough neighbors within distance ϵ ($\geq \text{minPts}$), it's a **core point** → start a new cluster.
 - If not, temporarily mark it as **noise** (it might later become a border point if another cluster grows into it).
- **Expand the cluster.**
 - Add all neighbors of the core point into a **queue of unvisited points**.
 - Each time we pull a point from the queue, we mark it *visited*.
 - If it's also a **core point**, add its neighbors to the queue (the cluster expands).
 - If it's a **border point**, just add it to the cluster, but don't expand further.
- **Repeat until queue is empty.**
 - Once there are no more unvisited points connected to the cluster, the cluster is complete.
- **Start again with another unvisited point.**
 - If it turns out to be a new core, grow another cluster.
 - If it's isolated, leave it as noise.

Key Intuition

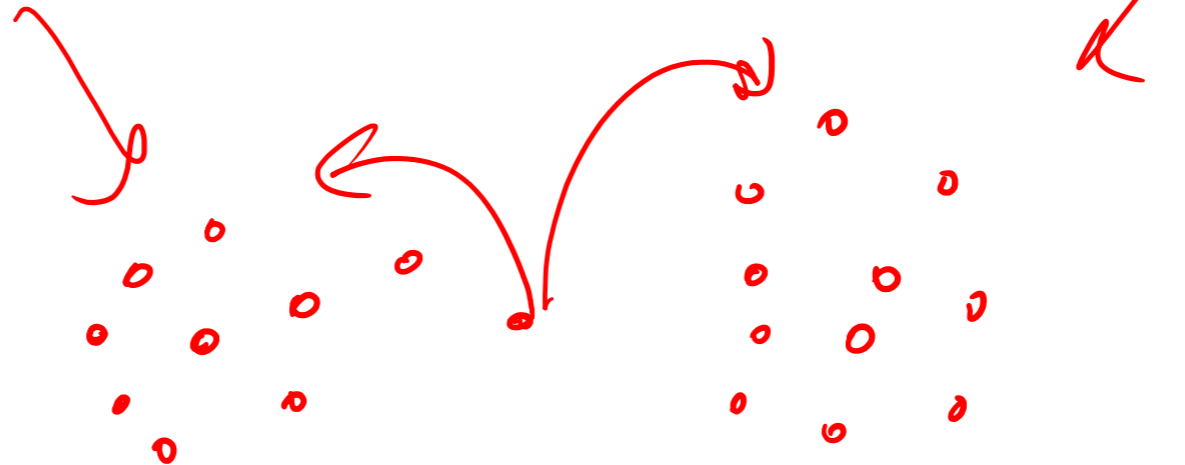
- **Core points** = dense regions (seeds of clusters).
- **Border points** = edges of clusters (belong to cluster but don't expand).
- **Noise** = isolated points (too few neighbors)

Deterministic?

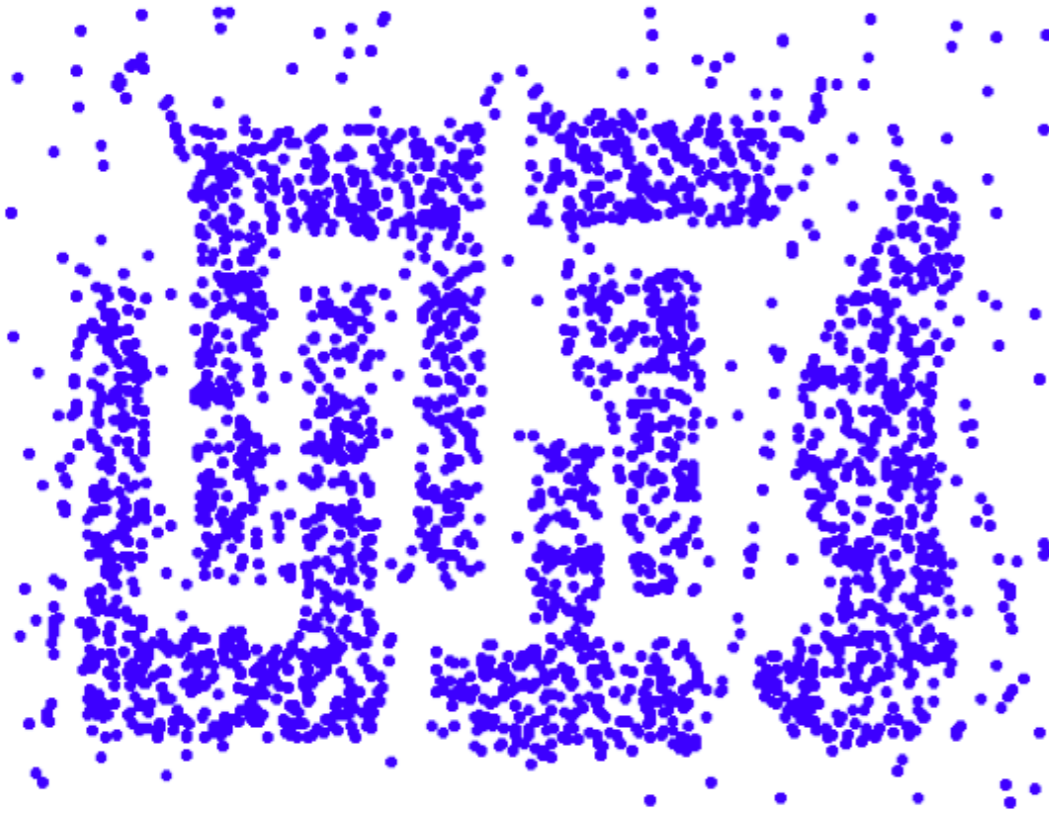
Same input \rightarrow same output



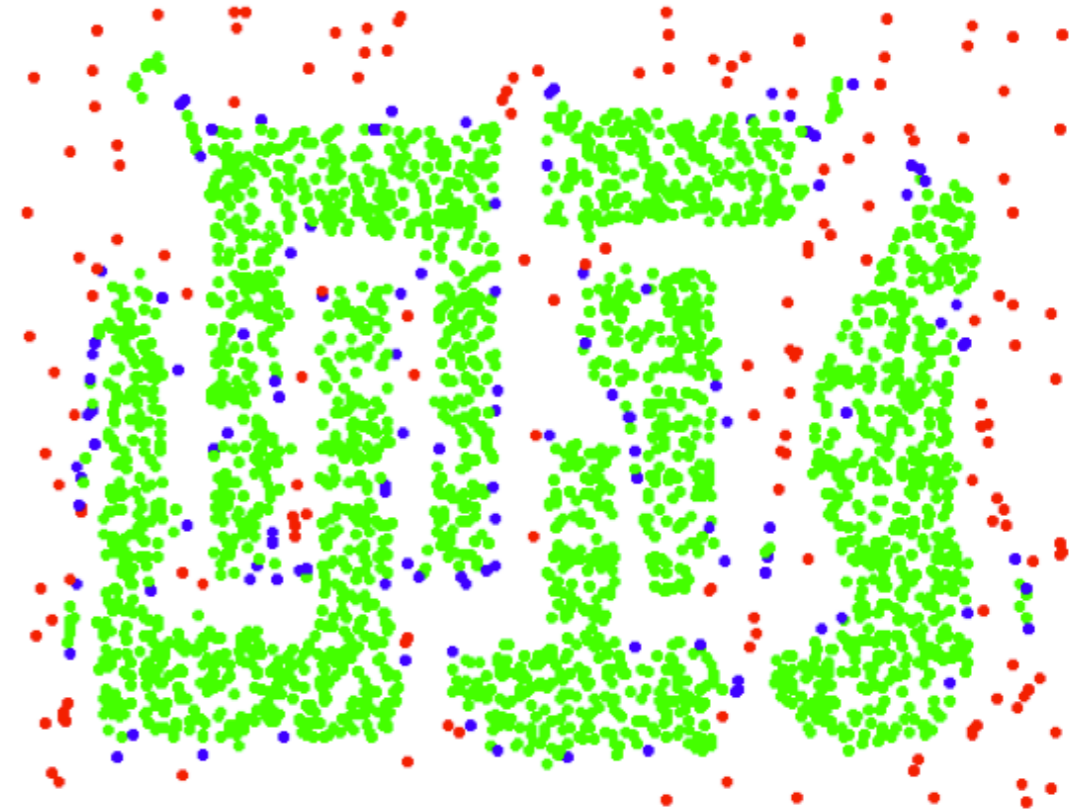
- **Almost deterministic**
- The result is deterministic **except for border points.**
- If a border point lies in the neighborhood of two clusters, it may get assigned differently depending on which cluster grows first.
- This is the only source of non-determinism.



Examples



Original Points



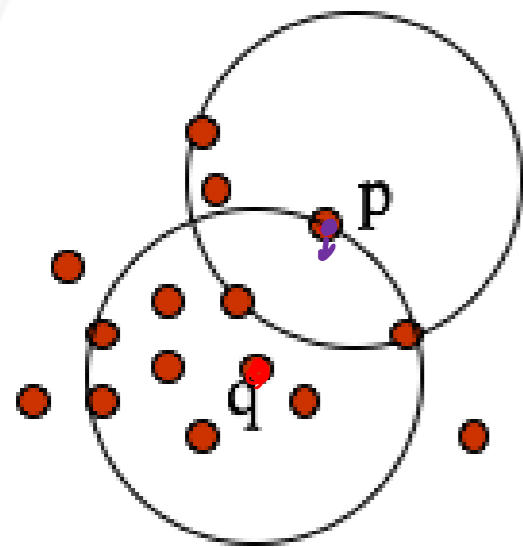
Point types: **core**,
border and **outliers**

$\epsilon = 10$, MinPts = 4

Density-based related points

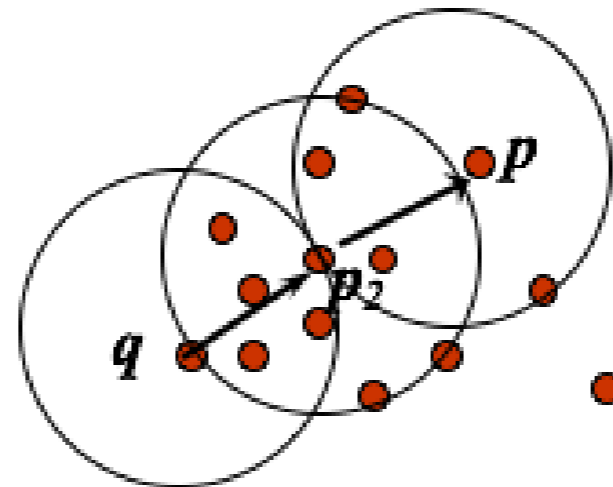
- Direct density reachability:
 - An object p is directly density-reachable from object q if (1) q is a core object; and (2) p is in q 's ϵ -neighborhood

p is direct density reachable from q

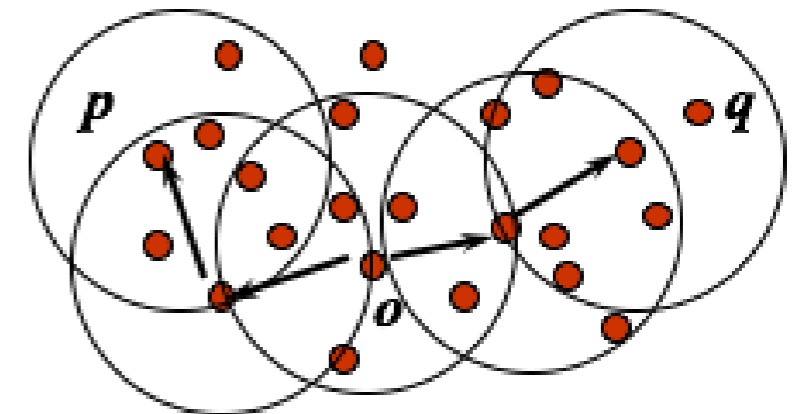


Directly Density-Reachable

MinPts = 5
Eps = 1 cm



Density-Reachable



Density-Connected

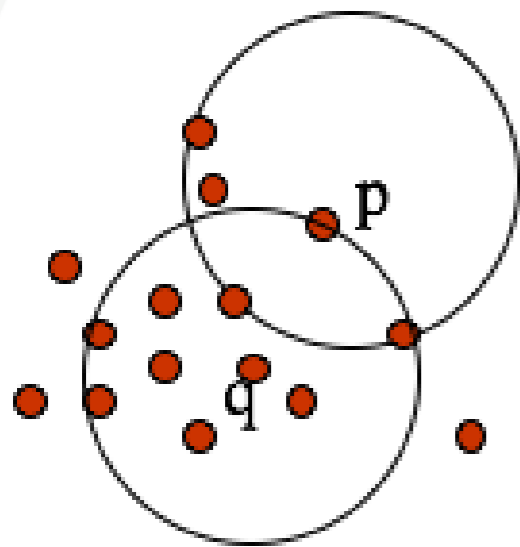
Density-based related points

- Density reachability:

- A point p is density-reachable from a point q if there is a chain of points $p_1, \dots, p_n, p_1 = q, p_n = p$ such that p_{i+1} is directly density-reachable from p_i

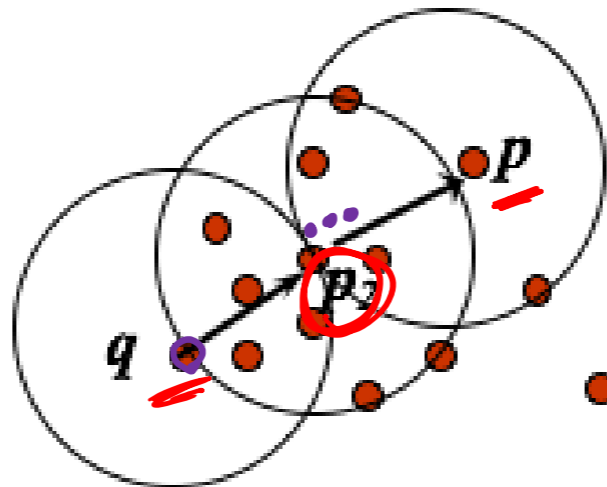
- $p_1 = q \rightarrow p_2 \rightarrow \dots \rightarrow p_n = p$

p is density reachable from q
~~X~~
Not a 2-way relationship

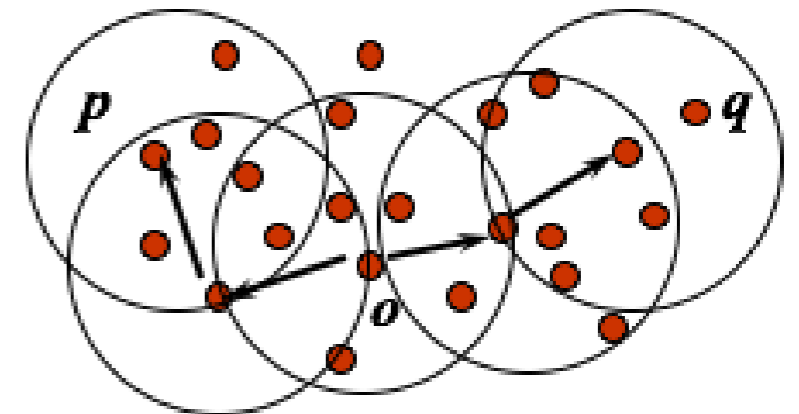


MinPts = 5
 Eps = 1 cm

Directly Density-Reachable



Density-Reachable



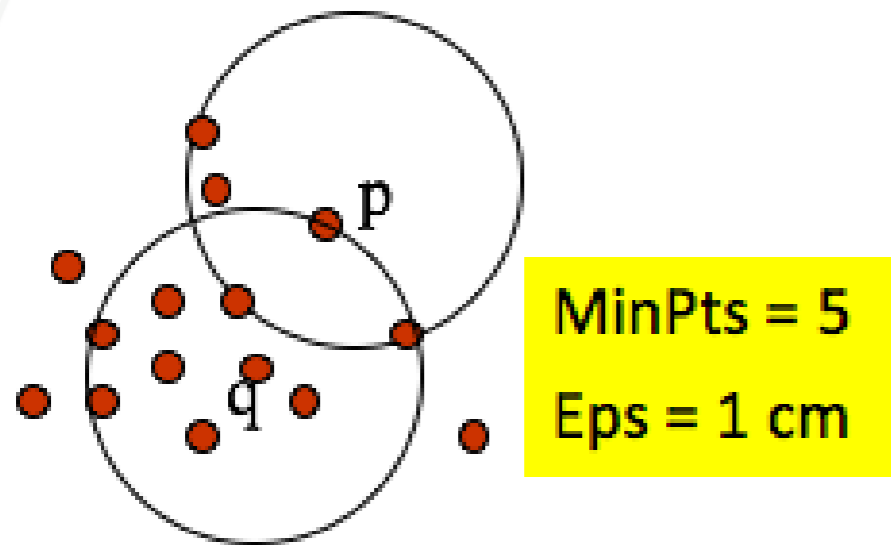
Density-Connected

Density-based related points

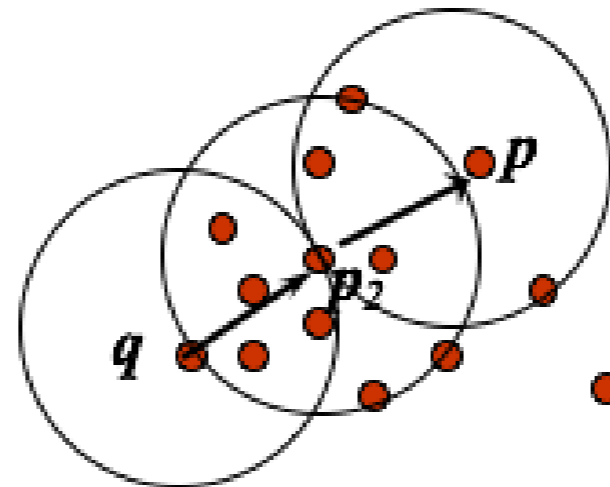
- Density connectivity:

- A point p is density-connected to a point q if there is a point o such that both p and q are density-reachable from o

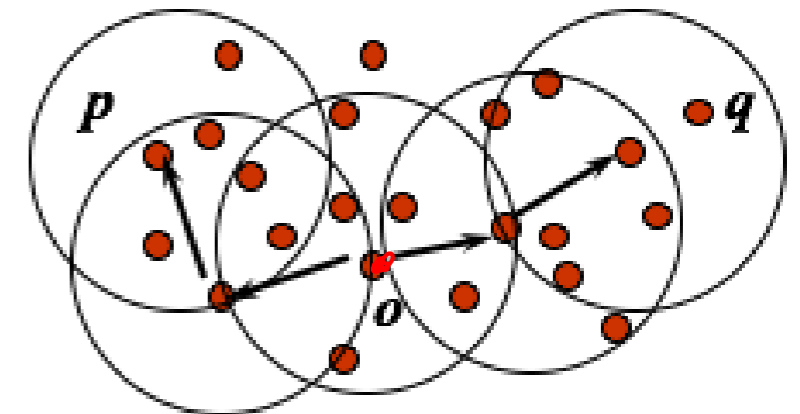
p & q are density connected



Directly Density-Reachable



Density-Reachable



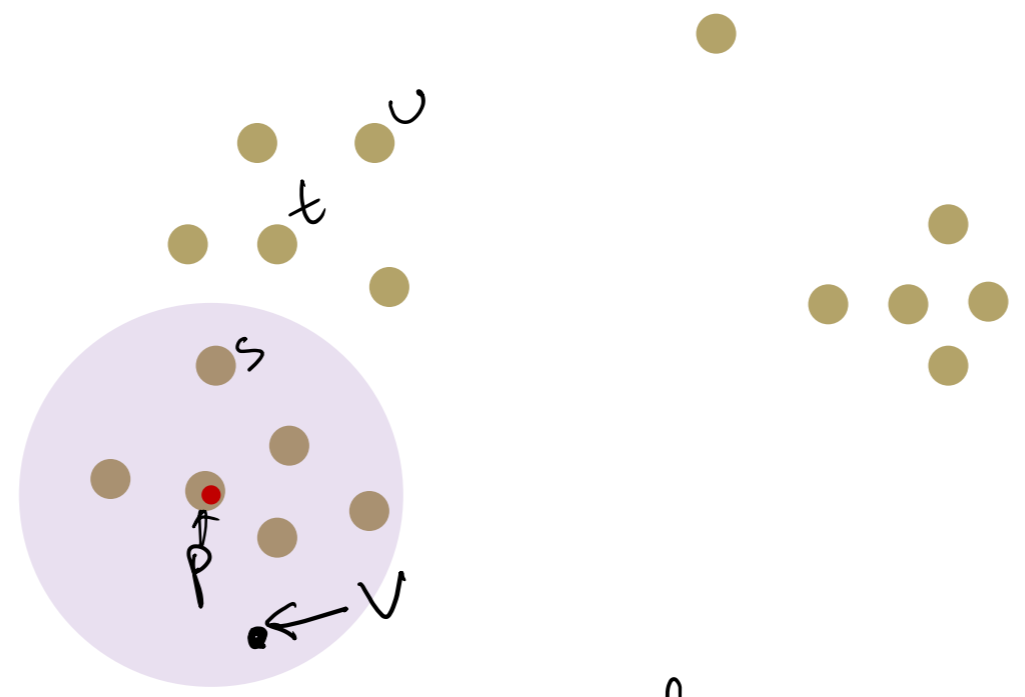
Density-Connected

t is direct density reachable from s

P is " " " " from s

U is " " " " from t
S is " " " " from t

$\epsilon = 1$ unit
MinPts = 5



U density reachable from s

v is ddx from p

v \xrightarrow{dx} s

v & u are density connected

Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN ←

The DBSCAN Algorithm

Algorithm is super sensitive to ϵ

DBSCAN(X , ϵ , MinPts)

$C = 0$

for each unvisited point P in dataset X

mark P as visited

NeighborPts = regionQuery(P , ϵ)

if sizeof(NeighborPts) < MinPts

mark P as NOISE

else

~~expandCluster(P , NeighborPts, C , ϵ , MinPts)~~

$C =$ next cluster

expandCluster(P , NeighborPts, C , ϵ , MinPts)

add P to cluster C

for each point P' in NeighborPts

if P' is not visited

mark P' as visited

NeighborPts' = regionQuery(P' , ϵ)

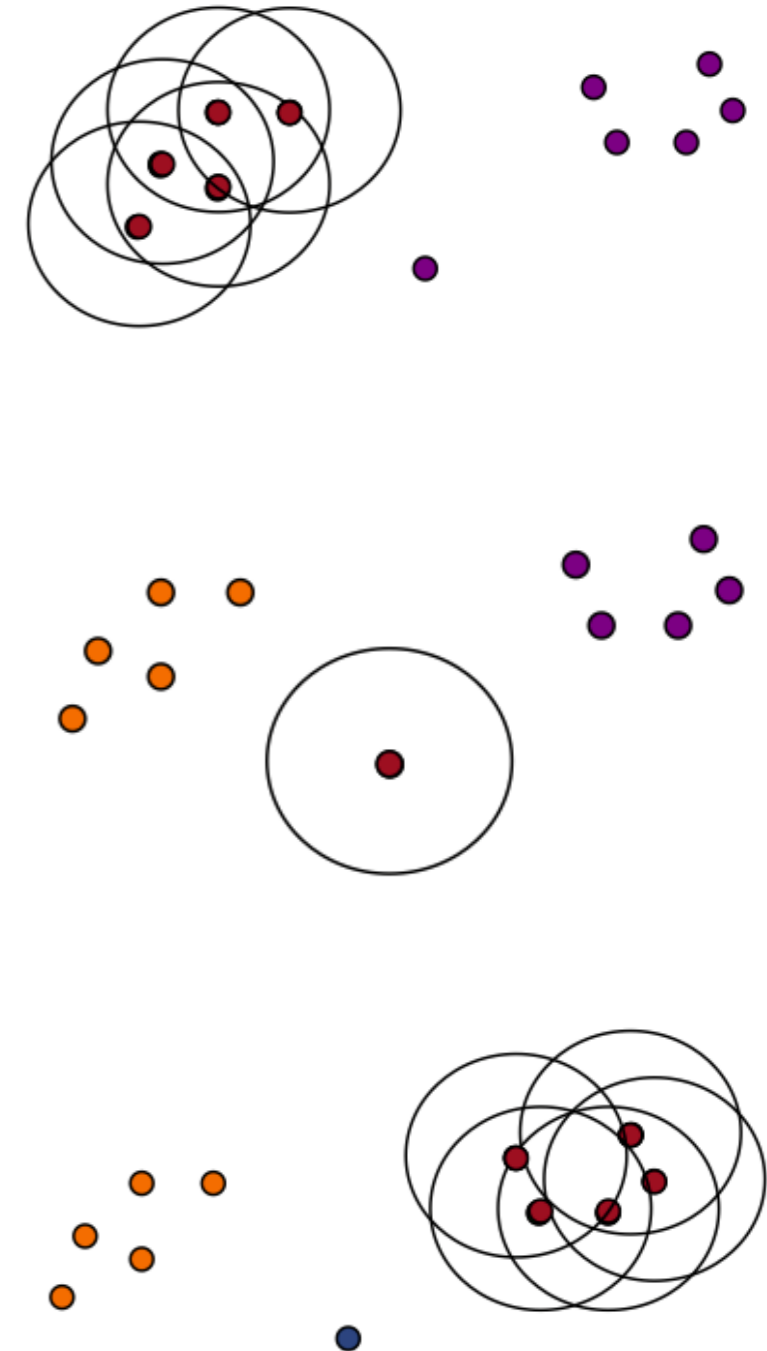
if sizeof(NeighborPts') \geq MinPts

NeighborPts = NeighborPts joined with NeighborPts'

if P' is not yet member of any cluster

add P' to cluster C

regionQuery(P , ϵ) return all points within P 's ϵ -neighborhood (including P)



which is very computationally expensive

Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN ←

DBSCAN is Sensitive to Parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

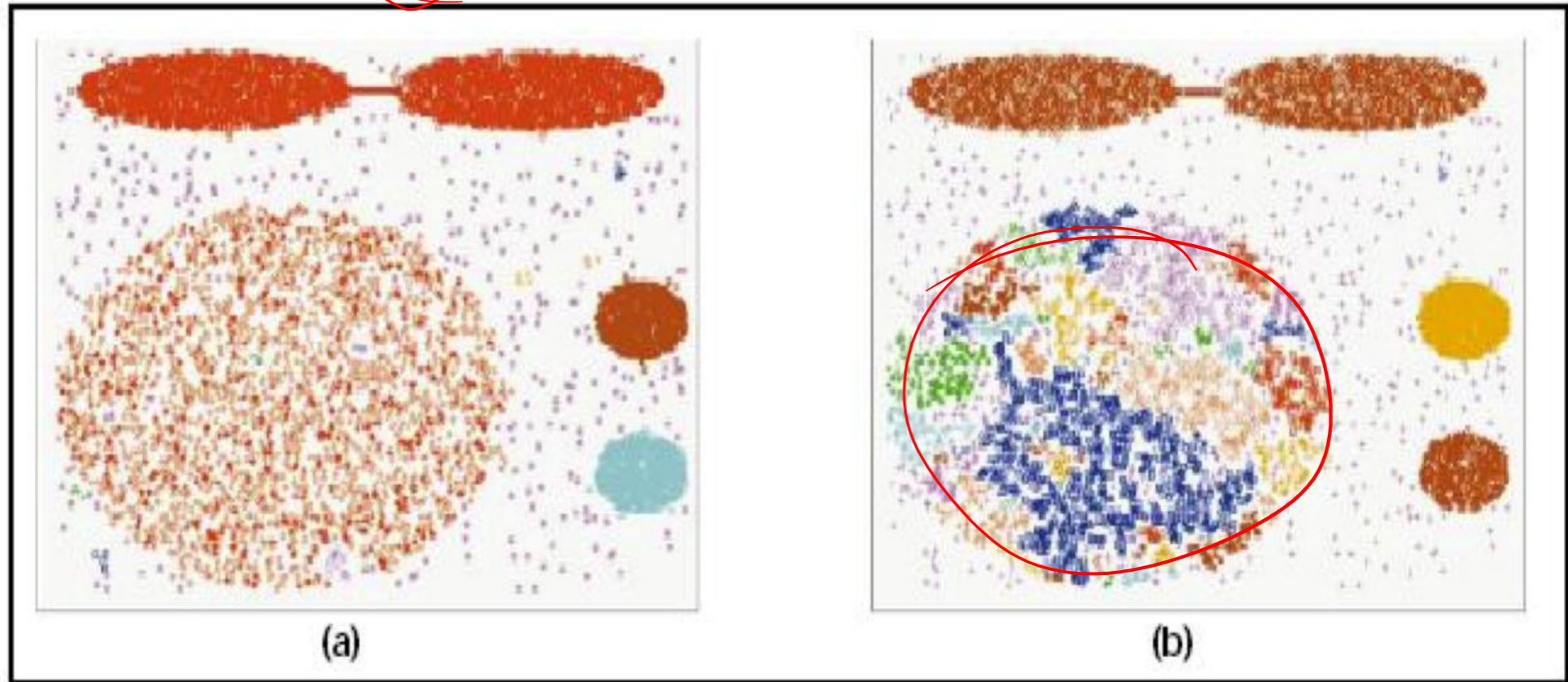
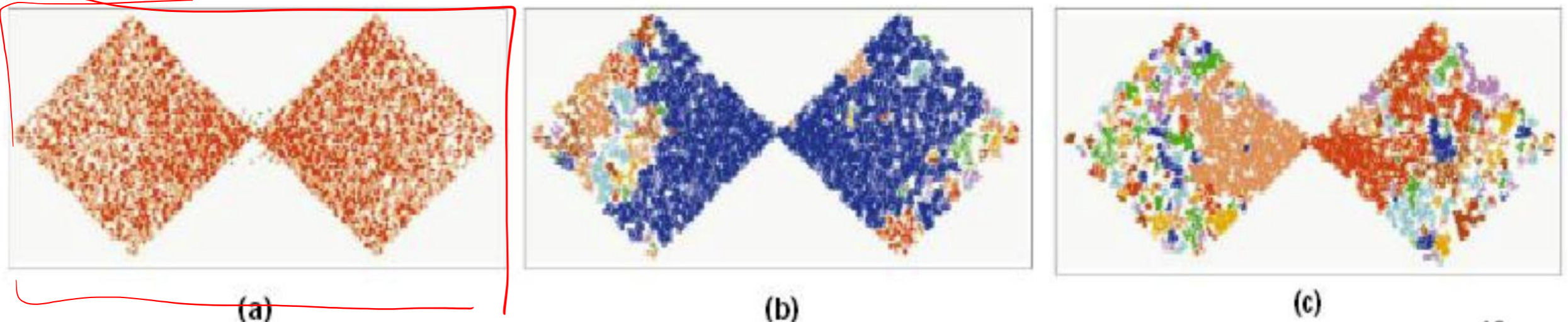
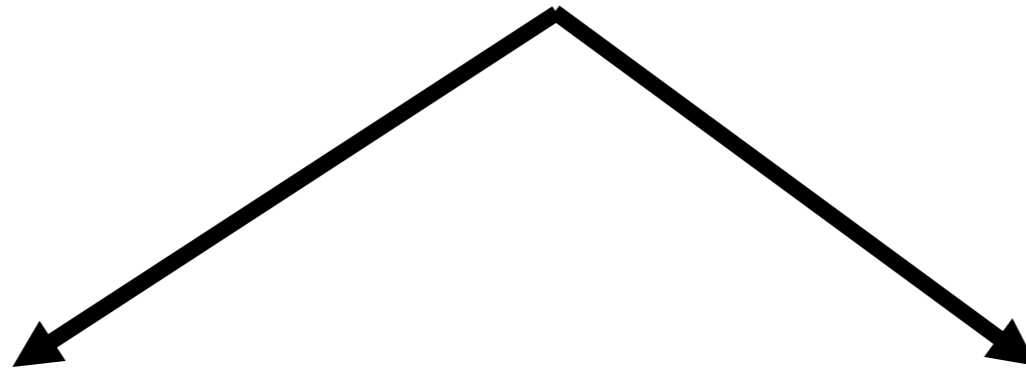


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



ϵ



High value (what will happen?)

Low value (what will happen?)

Clusters will merge and the majority of data points will be in the same cluster

A large part of data won't be clustered and considered as outliers. Because, they won't satisfy the number of points to create a dense region

Do we need to define the number of clusters in DBSCAN?

Nope

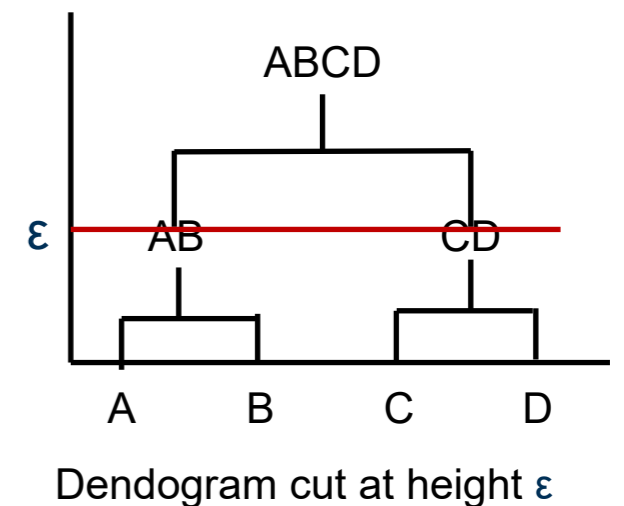
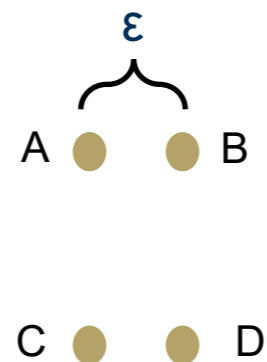
Minimum number of Points (**MinPts**)

Every point will be a cluster on its own, Why?

MinPts = 1?

Don't forget, in DBSCAN, a core point is counted as the number of neighboring points

MinPts = 2?



So, MinPts should be at least 3

no. of dimensions

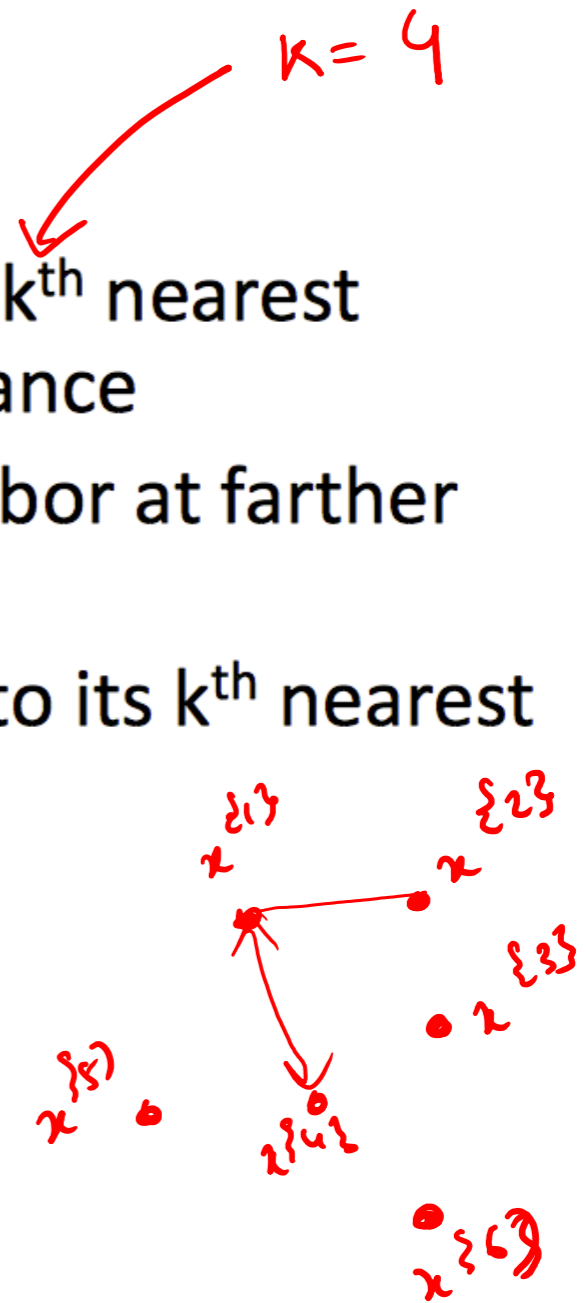
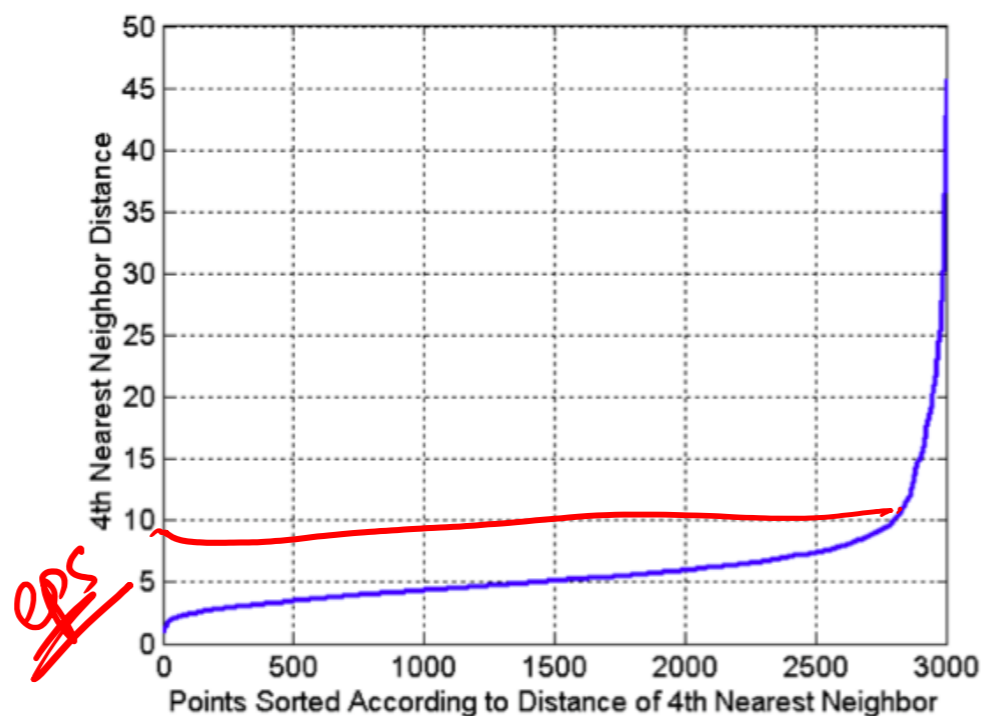
Rule of thumb, $\text{MinPts} \geq D+1$;
For noisy data \Rightarrow $\text{MinPts} = 2 \cdot D$ (yield more significant clusters)

How about Eps? (Elbow effect)

- Idea is that for points in a cluster, their k^{th} nearest neighbors are at roughly the same distance
- Noise points have the k^{th} nearest neighbor at farther distance
- So, plot sorted distance of every point to its k^{th} nearest neighbor

$x_{\{1\}}$
 $x_{\{2\}}$

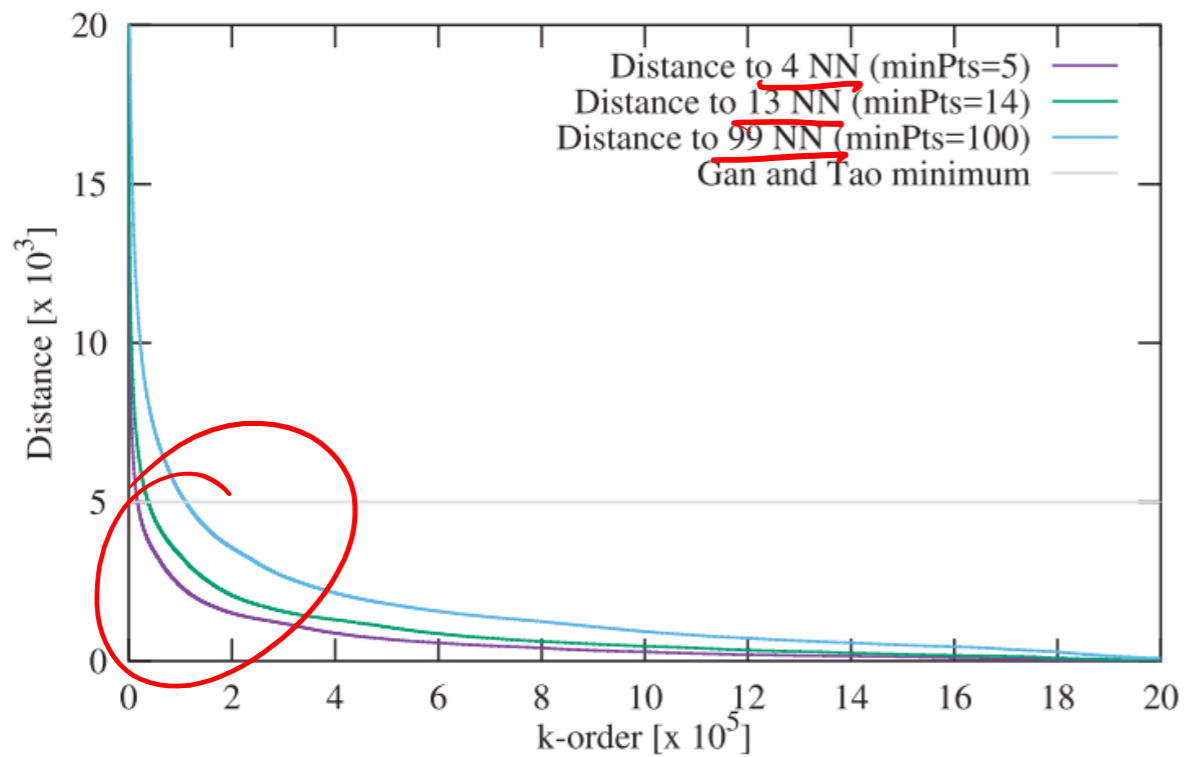
4th NN
distance



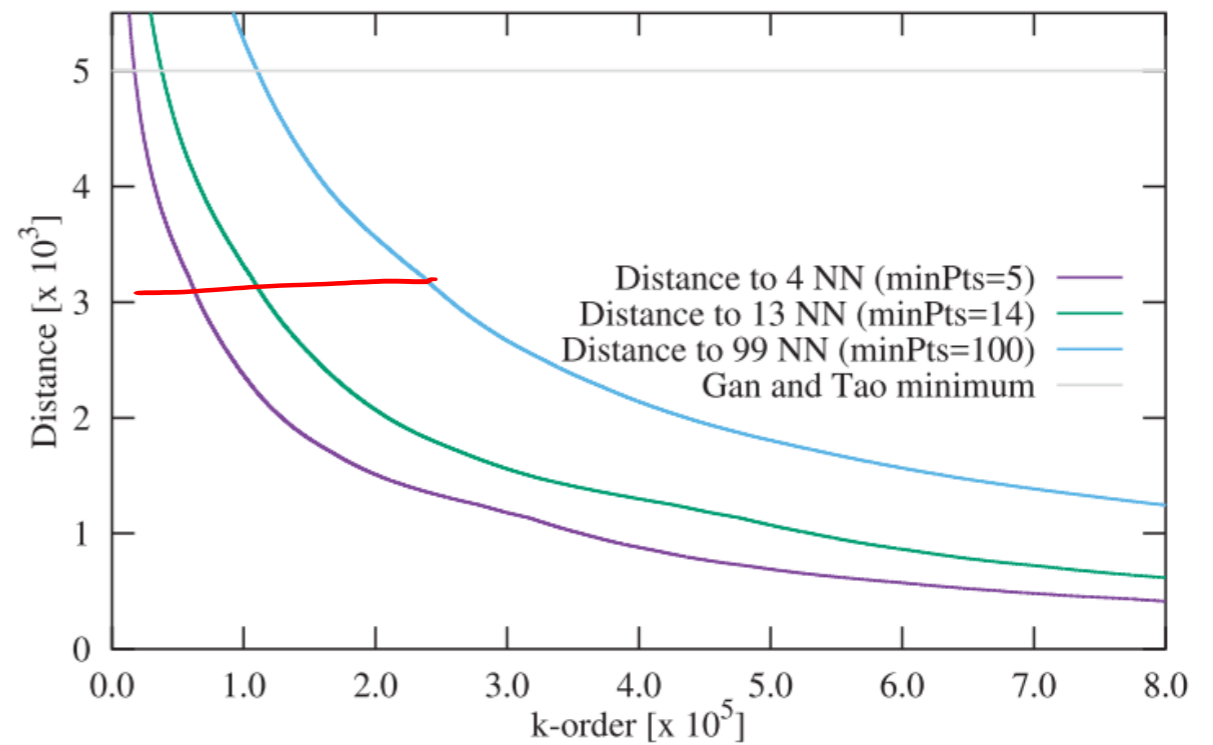
Here we have 3000 points and x-axis shows just a point index.

25 Point indices are sorted in ascending order based on their 4th nearest neighbor distance

Elbow effect another example



(a) k -distance plots



(b) k -distance plots (magnified region)

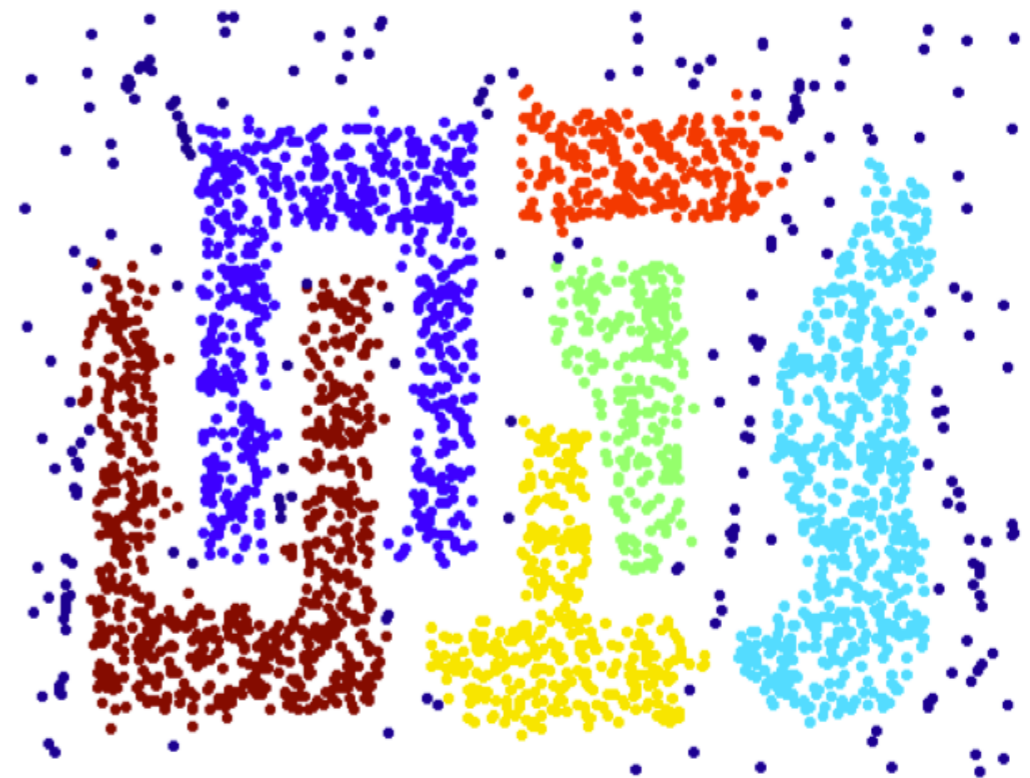
minPts often does not have a significant impact on the clustering results

When DBSCAN Works Well

- Robust to noise
- Can detect arbitrarily-shaped clusters



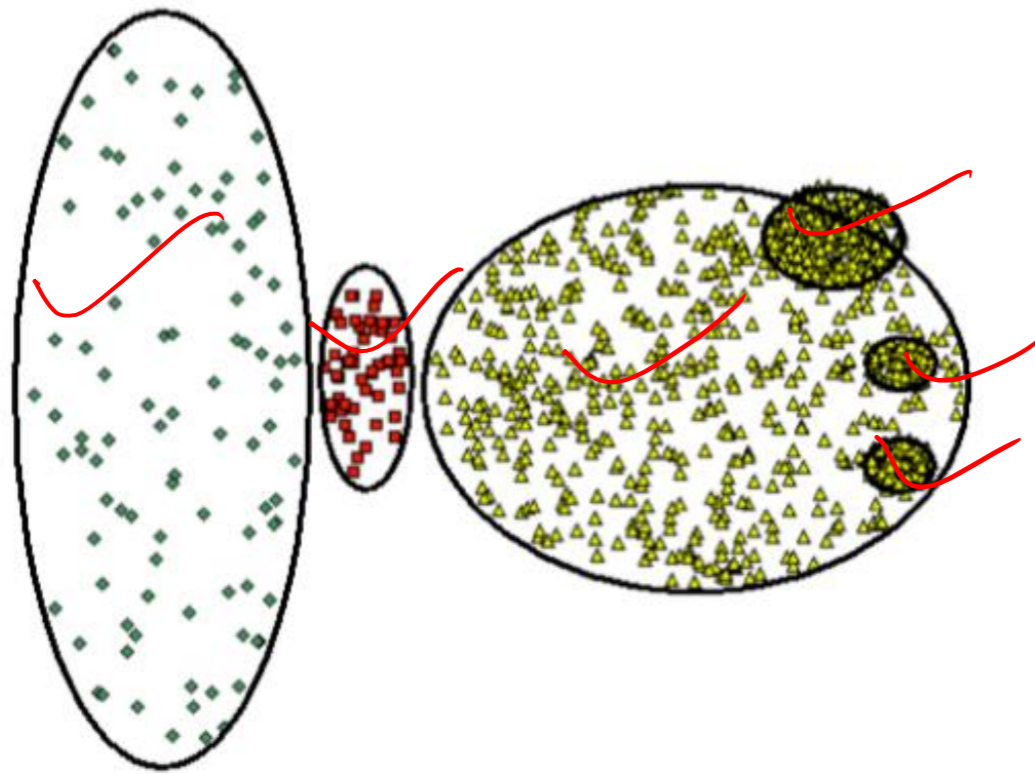
Original Points



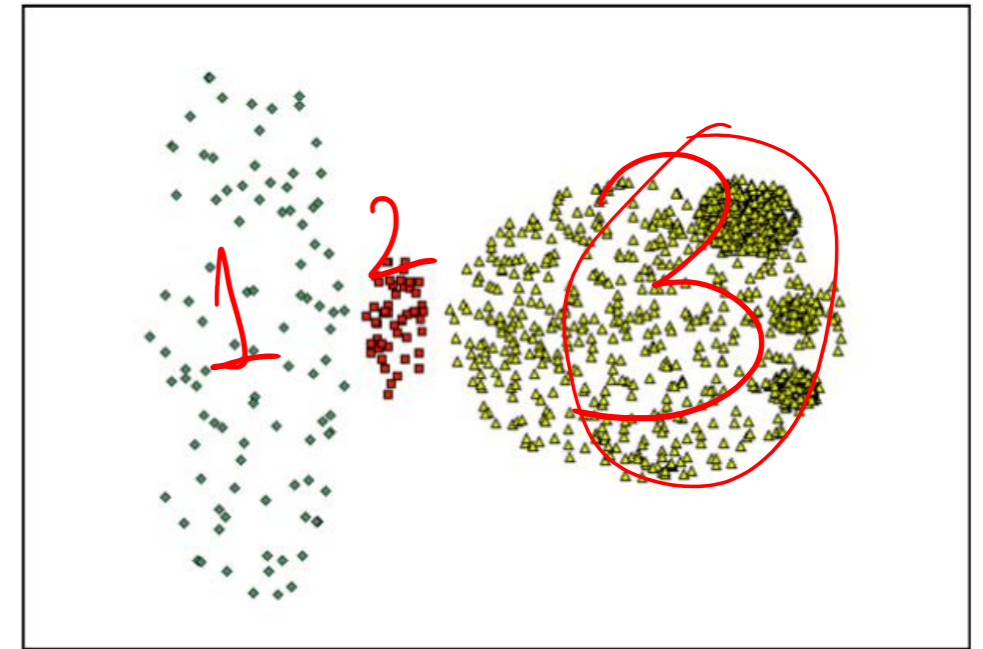
Clusters

When DBSCAN Does NOT Work Well

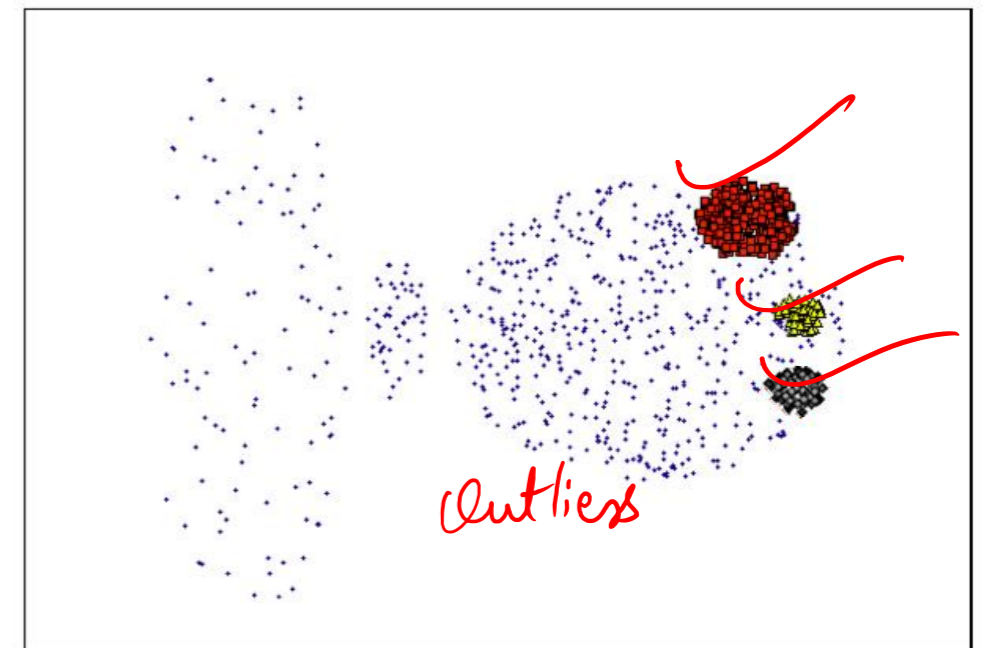
- Cannot handle varying densities
- Sensitive to parameters—hard to determine the best setting of parameters



Original Points



(MinPts=4, Eps=9.92).



(MinPts=4, Eps=9.75)

Take-Home Messages

- The basic idea of density-based clustering
- The two important parameters and the definitions of neighborhood and density in DBSCAN
- Core, border and outlier points
- DBSCAN algorithm
- DBSCAN's pros and cons



K-Means

"You guys always act like you're better than me"



GMM

DBSCAN

CHAMELEON

Which clustering method to choose?

Algorithm

Strengths

Best Contexts

Limitations

K-Means

Simple, fast, scalable to large datasets

Well-separated, spherical clusters; when you know/guess (k)

Struggles with non-spherical shapes, sensitive to outliers

GMM (Gaussian Mixture Models)

Probabilistic clusters, allows soft assignments

Data with overlapping clusters; when clusters are elliptical

Sensitive to initialization, assumes Gaussian distribution

Hierarchical Clustering

Produces dendrogram (multi-level clustering), no need to pre-specify (k)

Small–medium datasets; when hierarchy/taxonomy is useful (e.g. document organization)

Computationally expensive for large (n), sensitive to linkage choice

DBSCAN

Can find arbitrarily shaped clusters, robust to noise/outliers

Spatial/geometric data; clusters of varying shapes and sizes

Struggles with varying densities, parameter tuning (epsilon, minPts)

Clustering Evaluation

- Internal measures for clustering evaluation

- Elbow method ✓

k means, hierarchical

- Silhouette Coefficient ✓

$\frac{\text{Intra cluster distances}}{\text{Inter cluster distances}}$

- Graph-based measures (Beta-CV and Normalized cut)

- Davies-Bouldin Index

closer to 1

$\frac{\text{normalized cohesion}}{\text{normalized separation}}$

LOW

separation
Cohesion + Separation

HIGH

We want intra-cluster datapoints to be as close as possible to each other and inter-clusters to be as far as possible from each other

The Davies-Bouldin Index



Let μ_i denote the cluster mean

$$\mu_i = \frac{1}{n_i} \sum_{\mathbf{x}_j \in C_i} \mathbf{x}_j$$

Let σ_{μ_i} denote the dispersion or spread of the points around the cluster mean

Should be low \downarrow dispersion

$$\sigma_{\mu_i} = \sqrt{\frac{\sum_{\mathbf{x}_j \in C_i} \delta(\mathbf{x}_j, \mu_i)^2}{n_i}} = \sqrt{\text{var}(C_i)}$$

distance

The Davies–Bouldin measure for a pair of clusters C_i and C_j is defined as the ratio

Calculate the DB of i cluster from other clusters

$$DB_{ij} = \frac{\sigma_{\mu_i} + \sigma_{\mu_j}}{\delta(\mu_i, \mu_j)}$$

LOW HIGH

Worst Case

$$D_i = \max_{i \neq j} DB_{ij}$$

DB_{ij} measures how compact the clusters are compared to the distance between the cluster means. The Davies–Bouldin index is then defined as

$$DB = \frac{1}{k} \sum_{i=1}^k D_i$$

a lower value means that the clustering is better

Quick Knowledge Check

1. DBSCAN's key idea is to find clusters as...
A. Groups minimizing MSE around means
B. Dense regions separated by low-density gaps. C. Points on a grid with equal spacing
2. A **border** point is:
A. Not within ϵ of any point. B. Within ϵ of a core point but not itself core. C. The first point visited
3. Increase ϵ (holding MinPts fixed). What most likely happens?
A. More clusters, more noise. B. Fewer clusters; clusters may merge. C. No effect. D. All points become noise
4. Which statement is **true** about DBSCAN?
A. You must pre-specify the number of clusters k . B. It struggles with arbitrarily shaped clusters. C. It is robust to noise and finds non-spherical clusters. D. It cannot mark outliers
5. The **k-distance (elbow) plot** helps choose ϵ by:
A. Picking the global maximum of distances. B. Picking the knee of sorted k -th nearest neighbor distances. C. Minimizing average distance to centroids. D. Maximizing silhouette score
6. **Density-reachability** means p is reachable from q if:
A. There's a path of core points each within ϵ of the next. B. They share the same nearest centroid. C. Their Euclidean distance is minimum. D. They are in the same k -means cluster