

# Topics:

- Agglomerative clustering
- DBSCAN

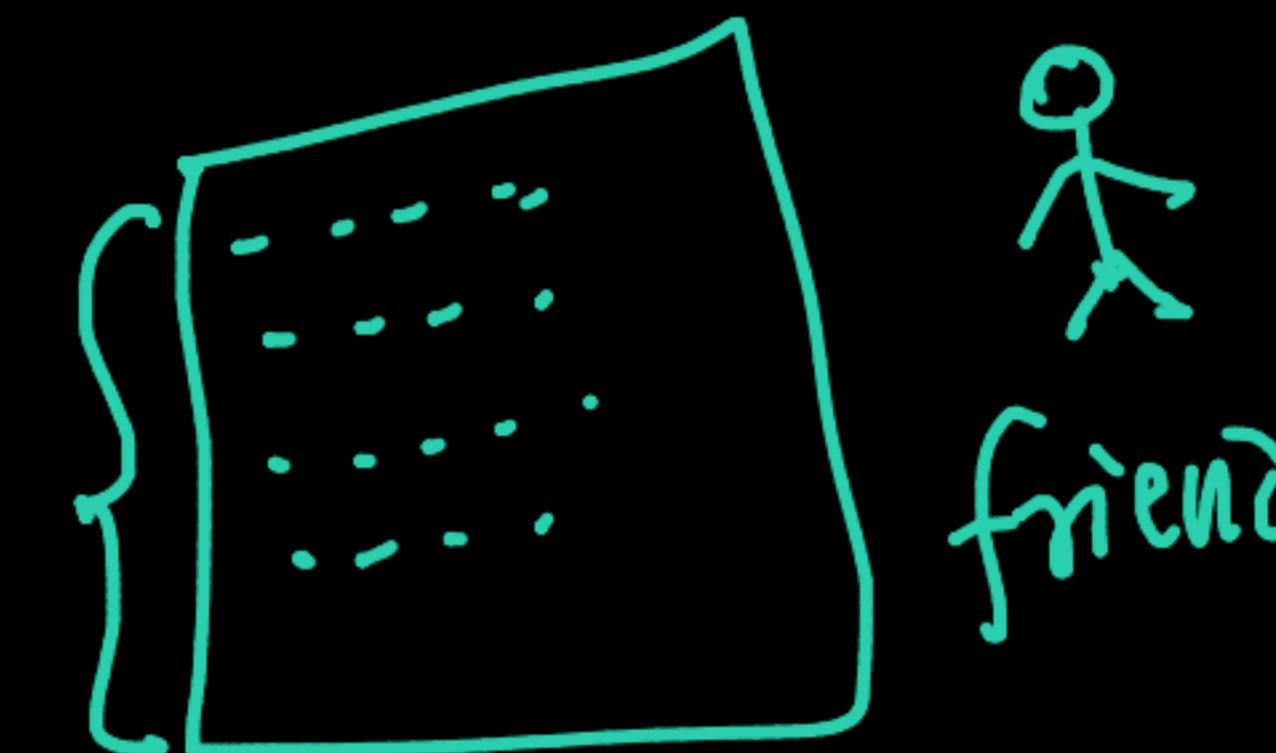


①

Spaced revision : →

3 weeks  
Week ;Month ; Quarter  
3M

②

Feynman's  
technique :-

friend

- first principles    "why"?    "how?"
- scribbled → video

Agglomerative



$n \rightarrow$

each point is a cluster

⋮

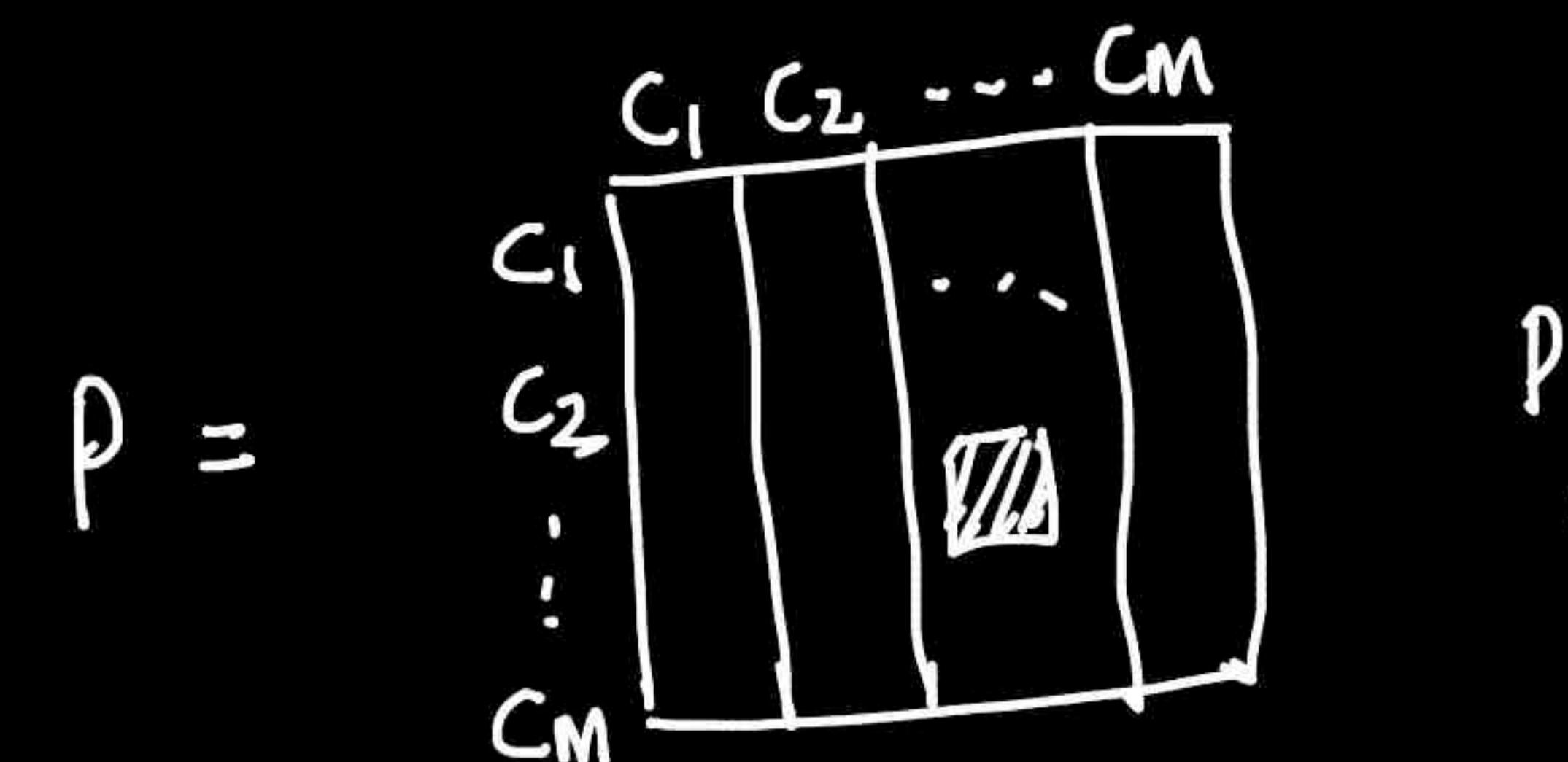
$1 \rightarrow$  one large cluster

## Basic - Structure:

{

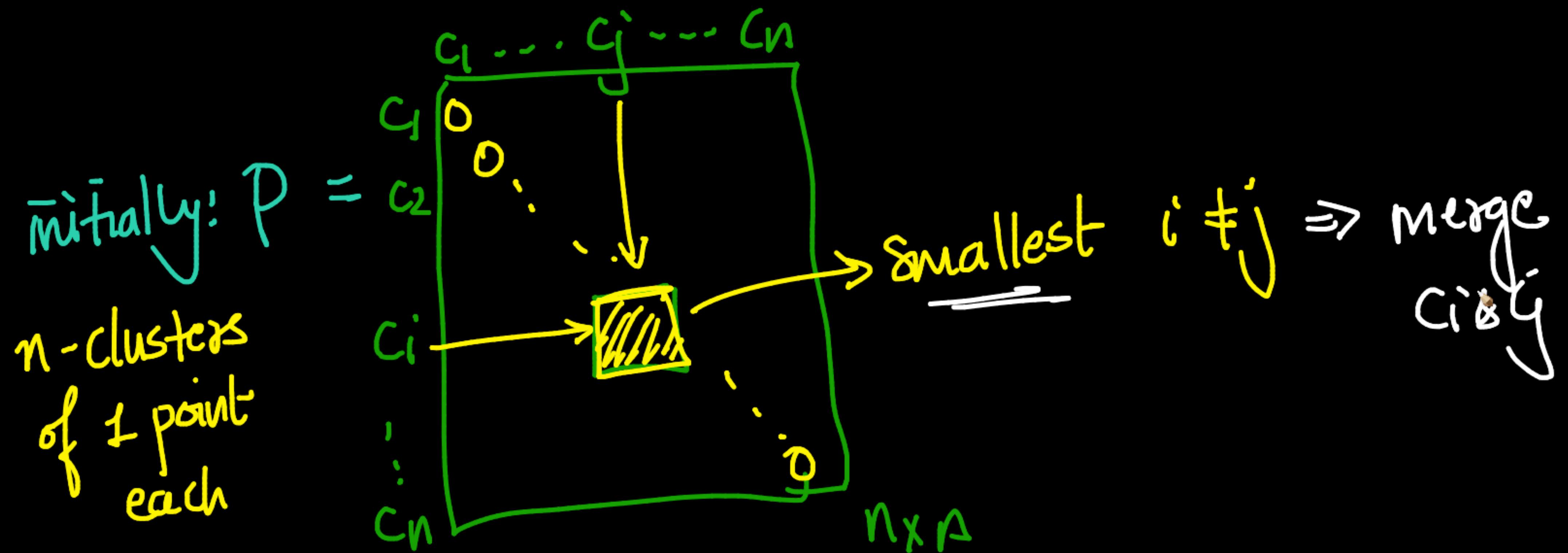
1. each point is a cluster :  $\cong$  data points
2. proximity-matrix  $P_{n \times n}$   $\cong$  (3) → how to compute distance b/w  $c_i \& c_j$
3. Repeat:
  - 3a. Merge the closest clusters
  - 3b. Update my proximity Matrix  $t$

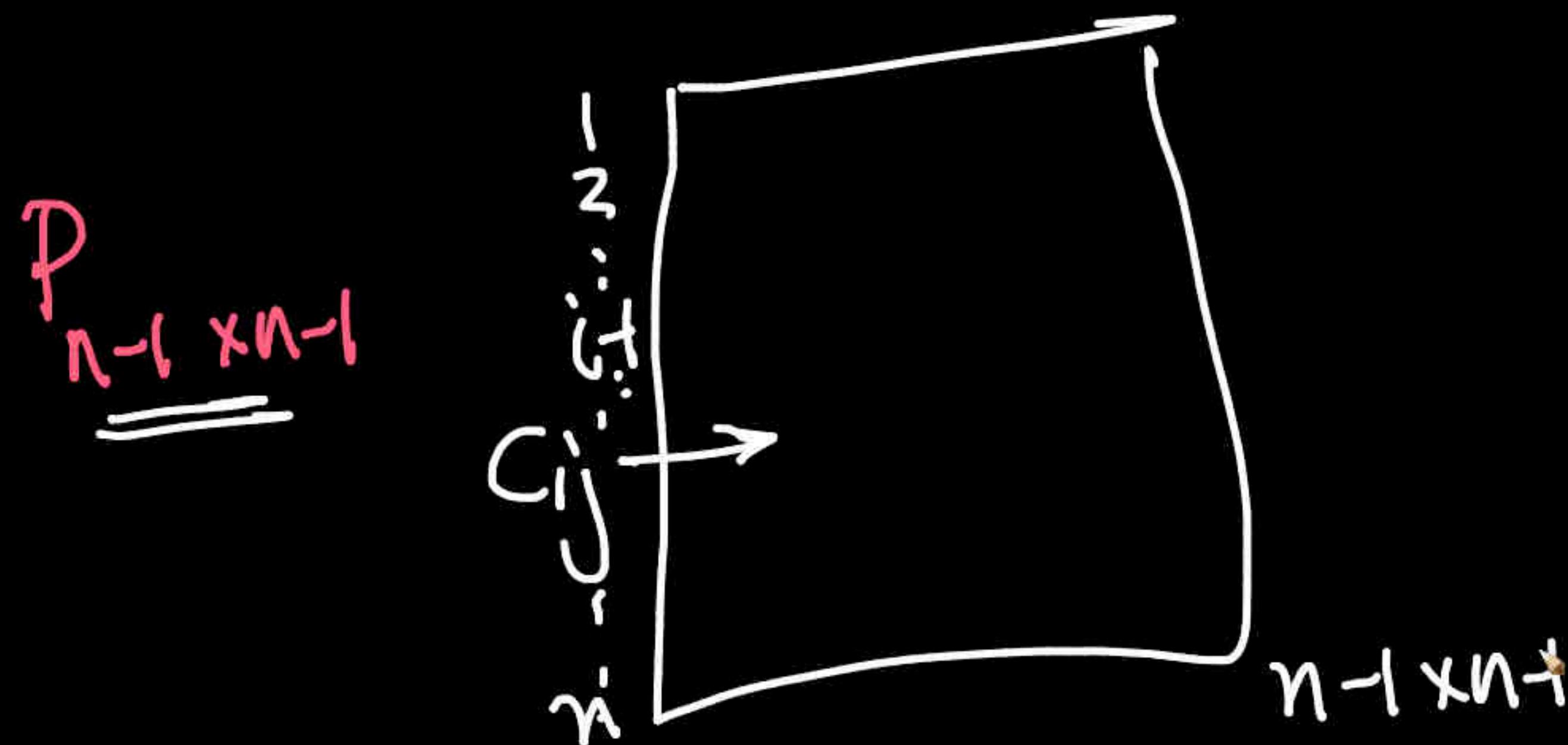
until a single cluster remains



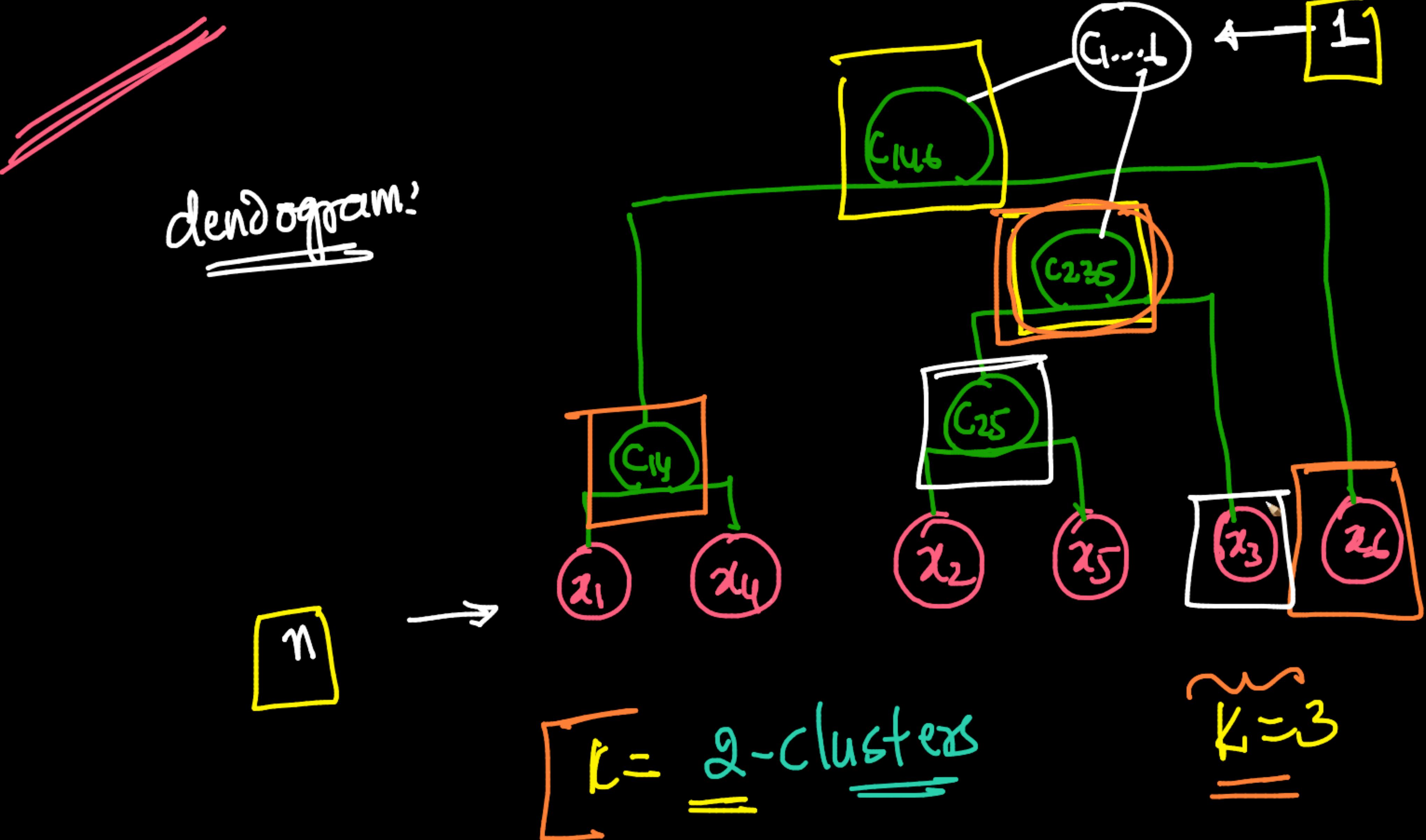
→ inv of distance

$$P_{ij} = \text{sim}(c_i, c_j)$$

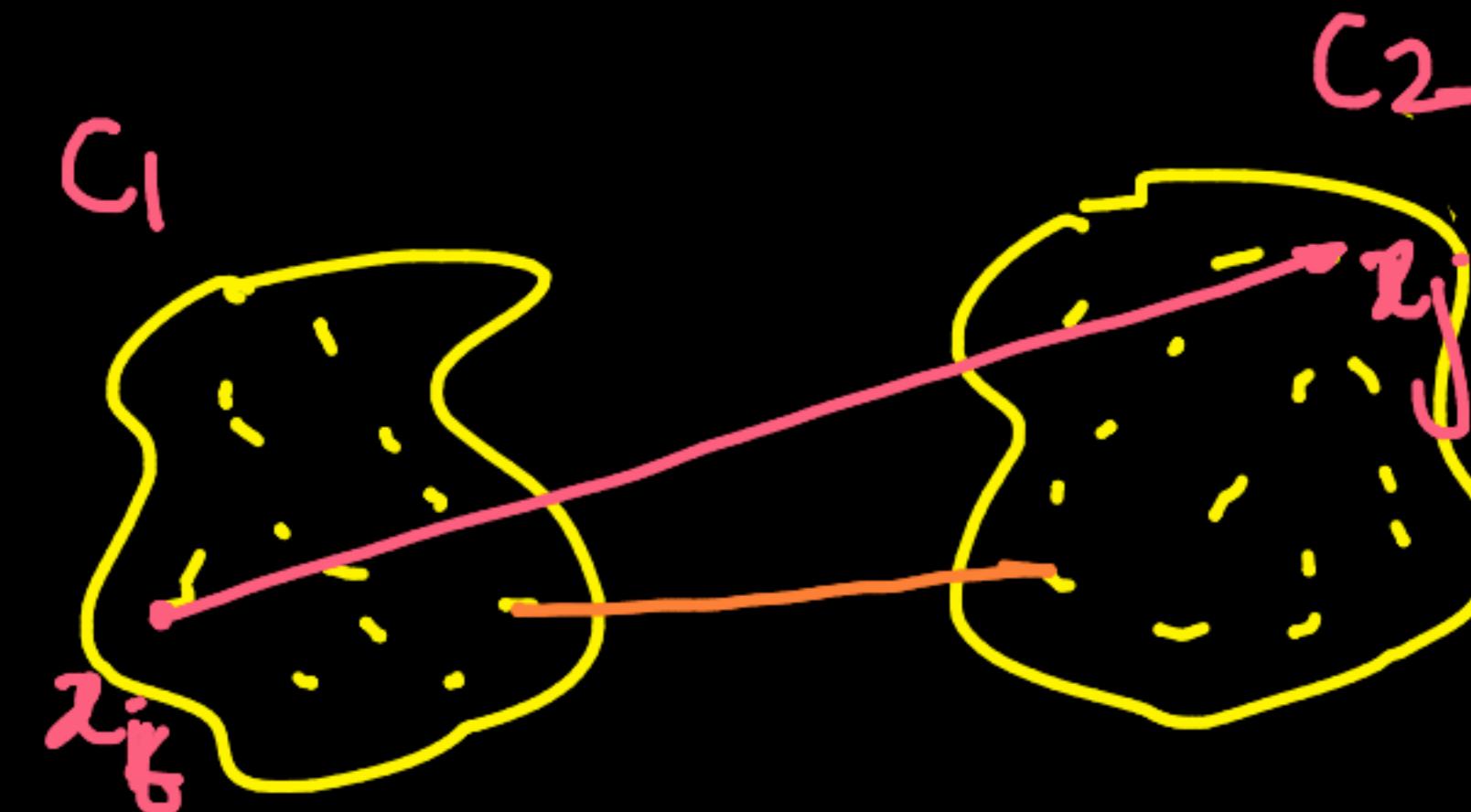




$$c_{ij} = \{x_i, x_j\}$$



(Q) How to compute  $\text{dist}(C_1, C_2)$ ?



5: Ward's dist euc. 2

$$\sum_{x_j \in C_2} \sum_{x_i \in C_1} \frac{\text{dist}(x_i, x_j)}{|C_1| |C_2|}$$

**3. MIN**:  $\min_{x_i \in C_1; x_j \in C_2} \text{dist}(x_i, x_j)$

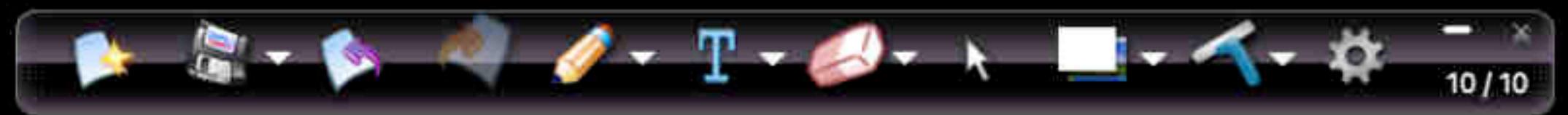
**4. group average dist**:  $\frac{\sum_{x_j \in C_2} \sum_{x_i \in C_1} \text{dist}(x_i, x_j)}{|C_1||C_2|}$

**1. euc. dist b/w centroids**:  $\sqrt{\sum_{i=1}^n (x_{i1} - x_{i2})^2}$

**2. MAX**:  $\max_{x_i \in C_1; x_j \in C_2} \text{dist}(x_i, x_j)$

**3. MIN**:  $\min_{x_i \in C_1; x_j \in C_2} \text{dist}(x_i, x_j)$

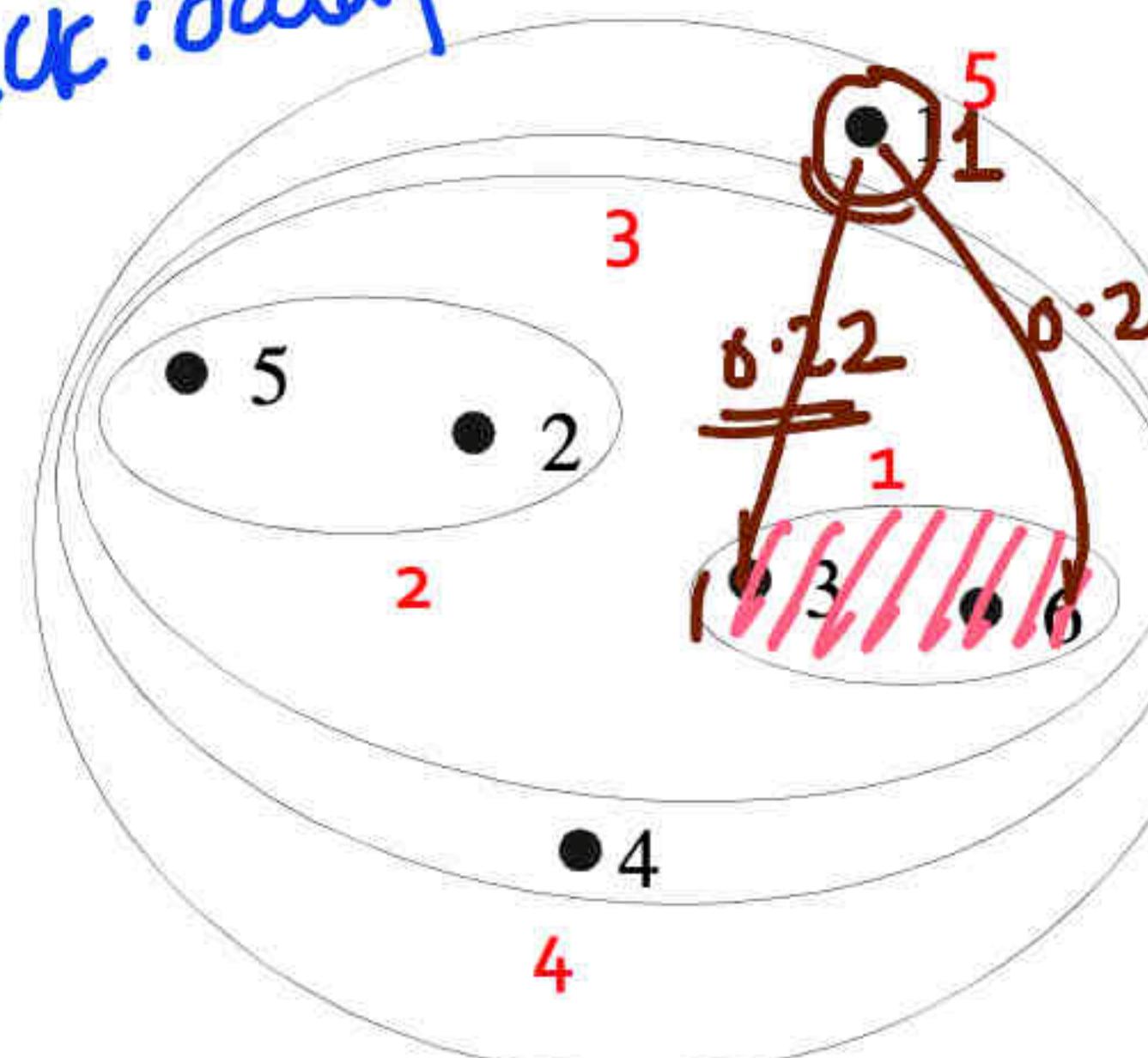
**4. group average dist**:  $\frac{\sum_{x_j \in C_2} \sum_{x_i \in C_1} \text{dist}(x_i, x_j)}{|C_1||C_2|}$



Simple & low-math  
(logical & intuitive)

# Hierarchical Clustering: MIN

black: data points



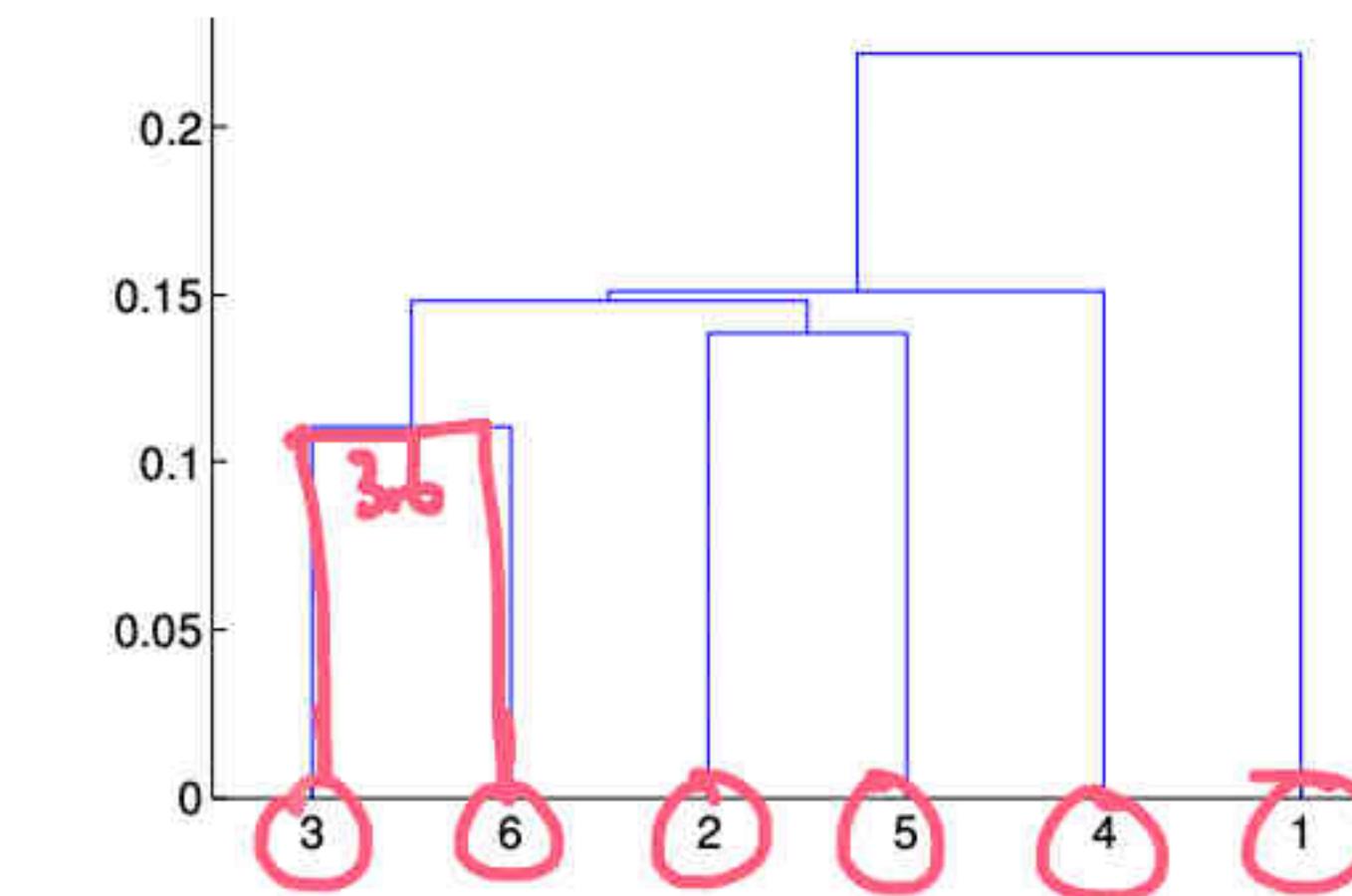
Nested Clusters

Dendrogram



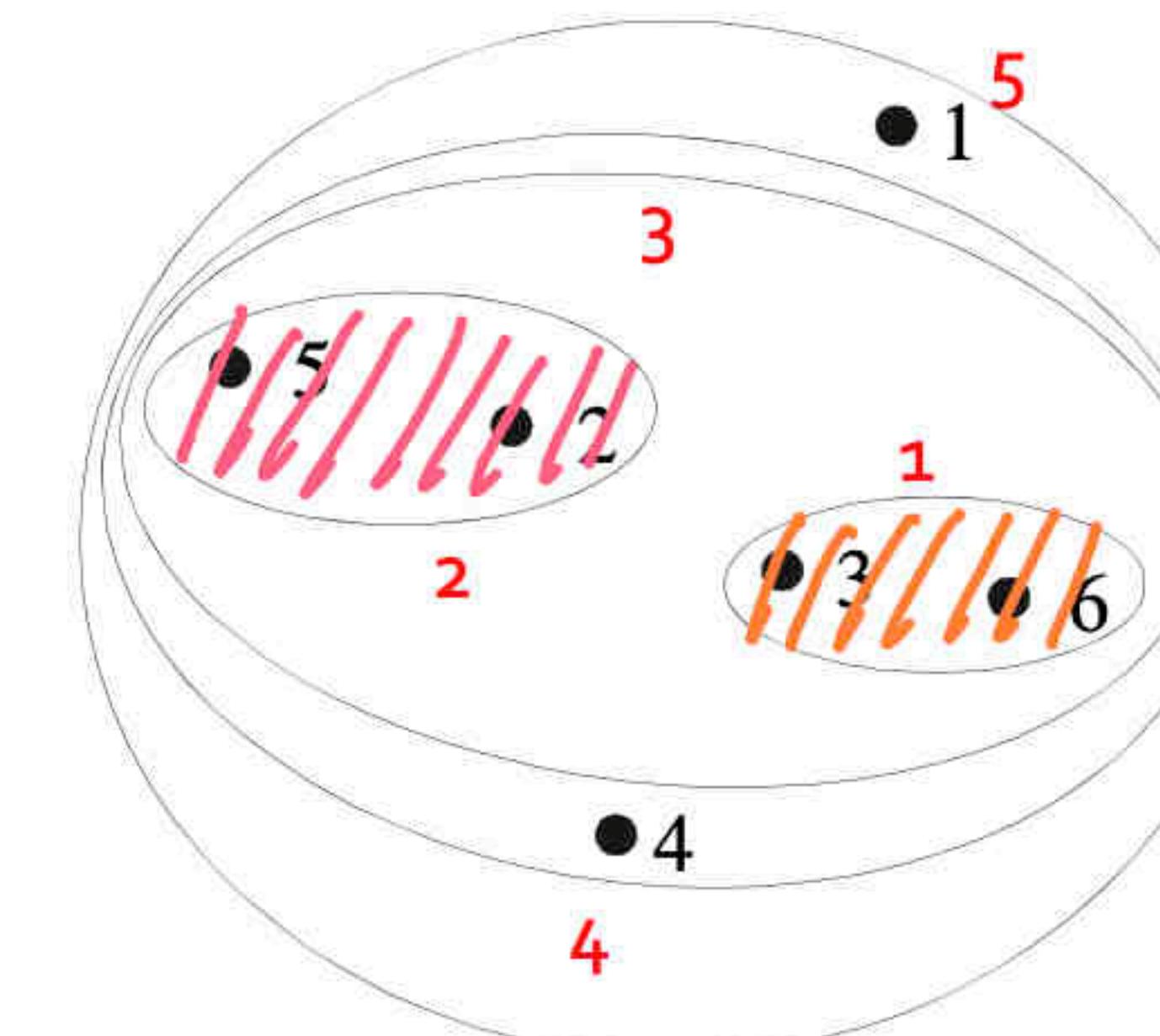
	1	2	3	4	5	6
1	0	.24	.22	.37	.34	.23
2	.24	0	.15	.20	.14	.25
3	.22	.15	0	.15	.28	.17
4	.37	.20	.15	0	.29	.22
5	.34	.14	.28	.29	0	.39
6	.23	.25	.17	.22	.39	0

5x5



# Hierarchical Clustering: MIN

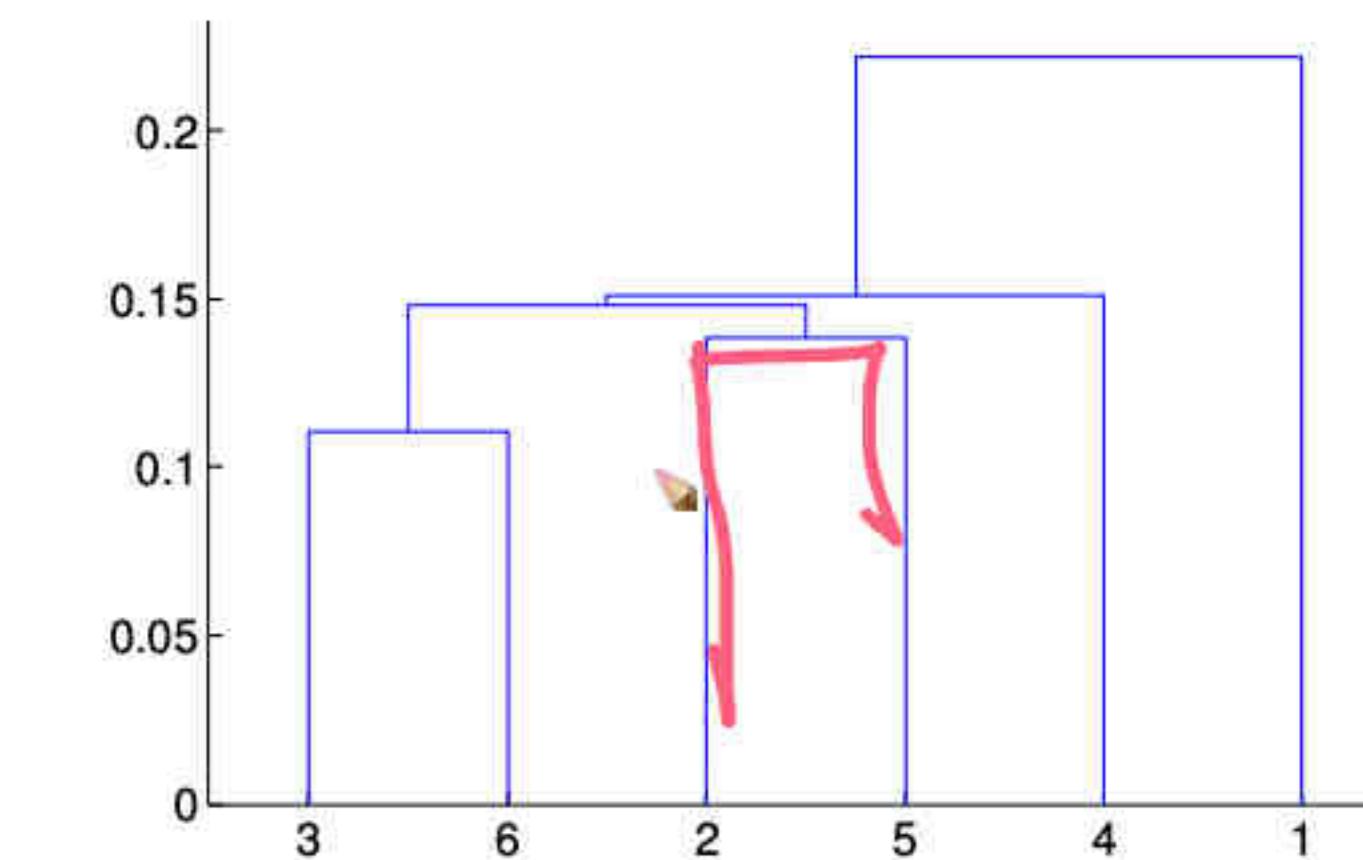
let's  
0.15  
 $C_1 - C_2$   
 $C_2 - C_3$   
0.15



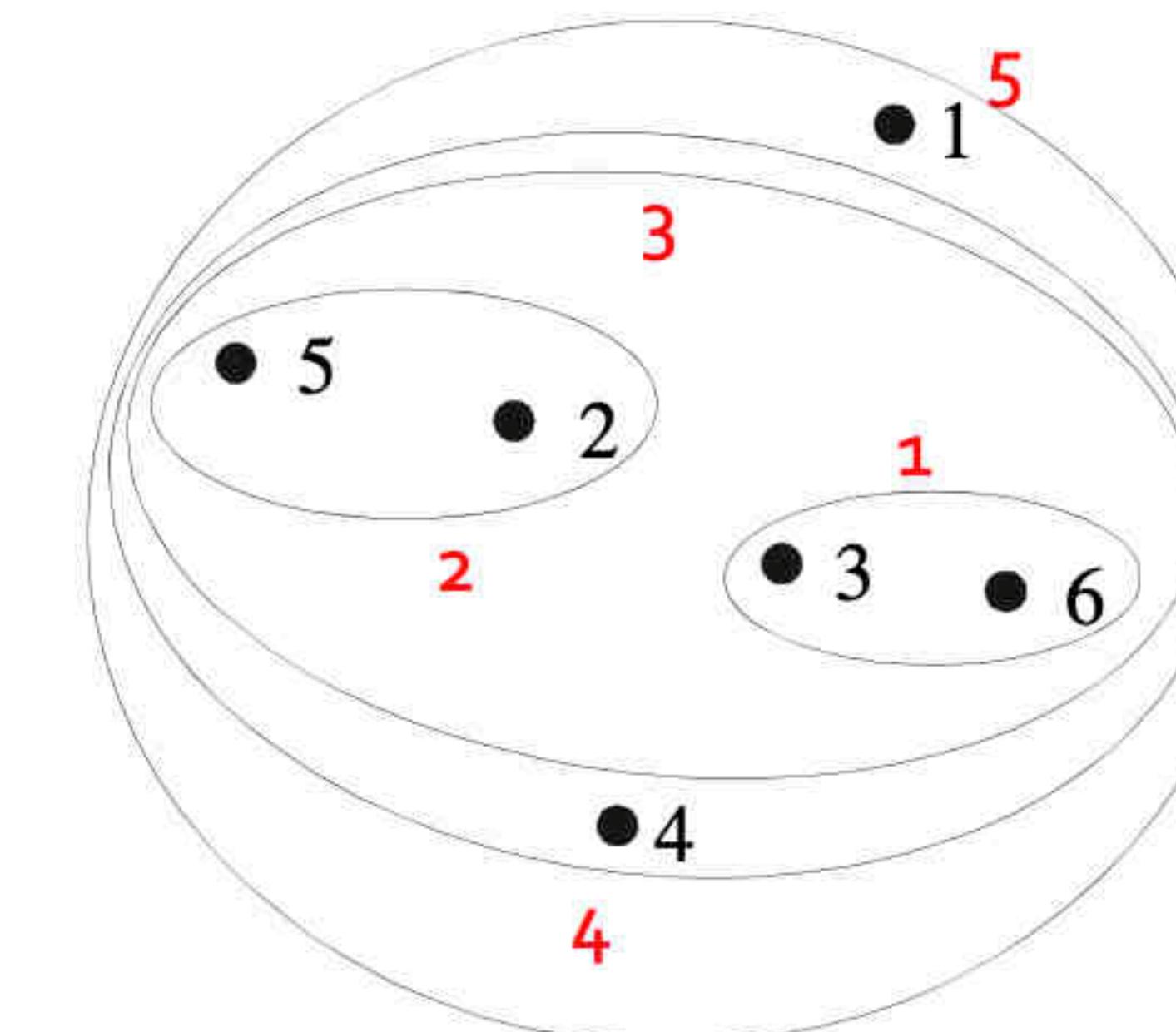
Nested Clusters

Dendrogram

	1	2	3	4	5	6
1	0	.24	.22	.37	.34	.23
2	.24	0	.15	.20	.14	.25
3	.22	.15	0	.15	.28	.21
4	.37	.20	.15	0	.29	.22
5	.34	.14	.28	.29	0	.39
6	.23	.25	.11	.22	.39	0



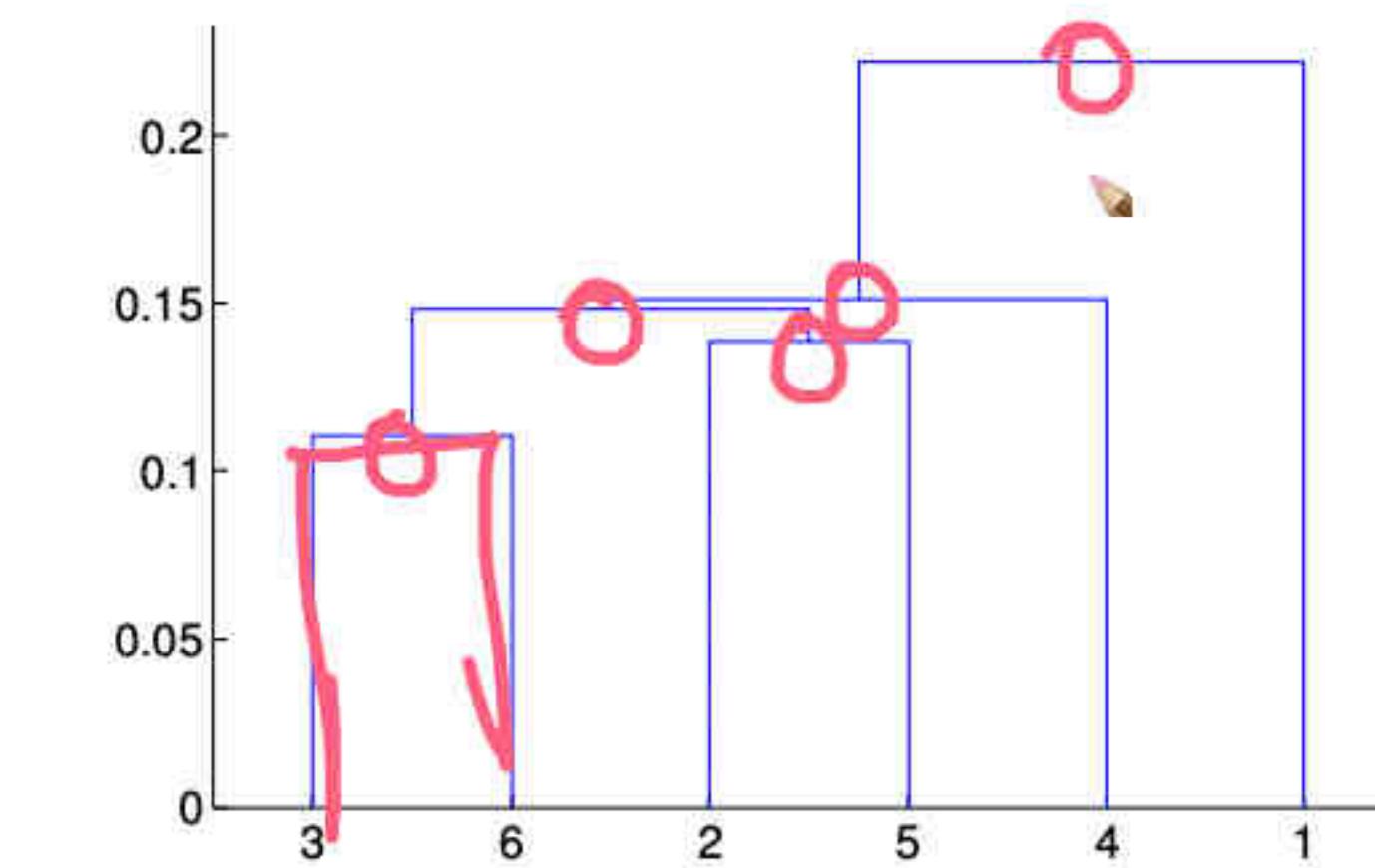
# Hierarchical Clustering: MIN



Nested Clusters

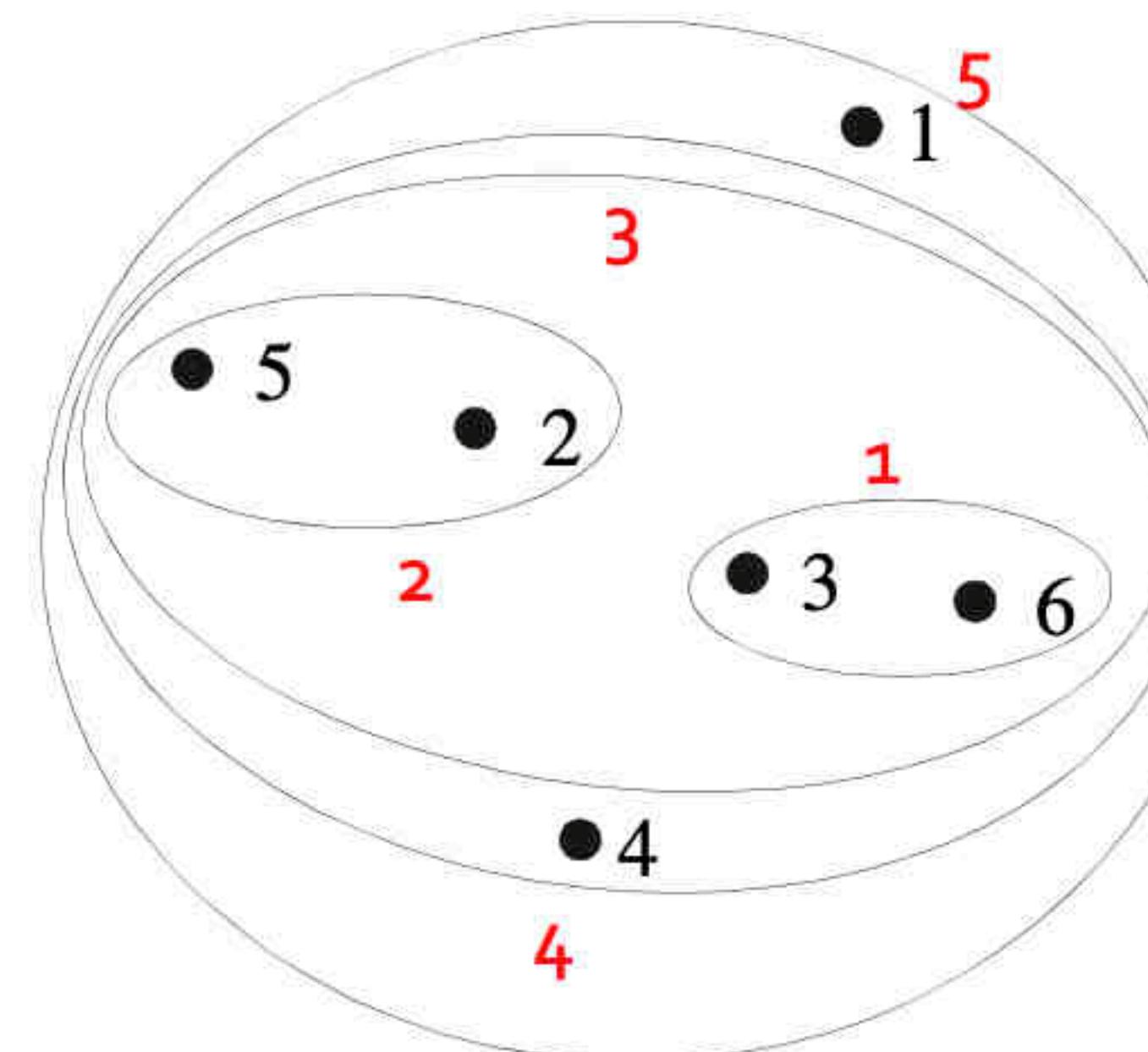
Dendrogram

	1	2	3	4	5	6
1	0	.24	.22	.37	.34	.23
2	.24	0	.15	.20	.14	.25
3	.22	.15	0	.15	.28	.11
4	.37	.20	.15	0	.29	.22
5	.34	.14	.28	.29	0	.39
6	.23	.25	.11	.22	.39	0



# Hierarchical Clustering: MIN

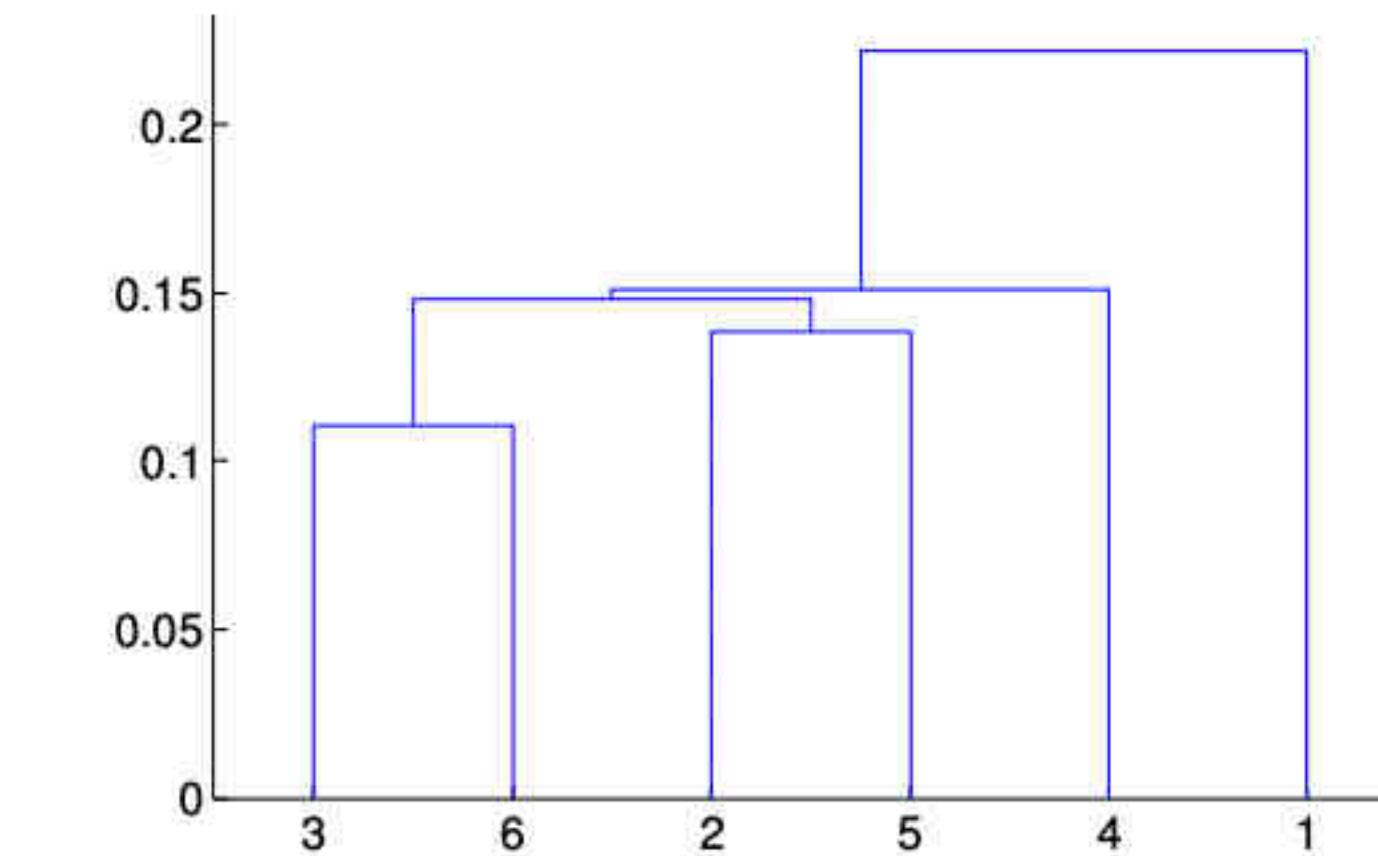
Manually



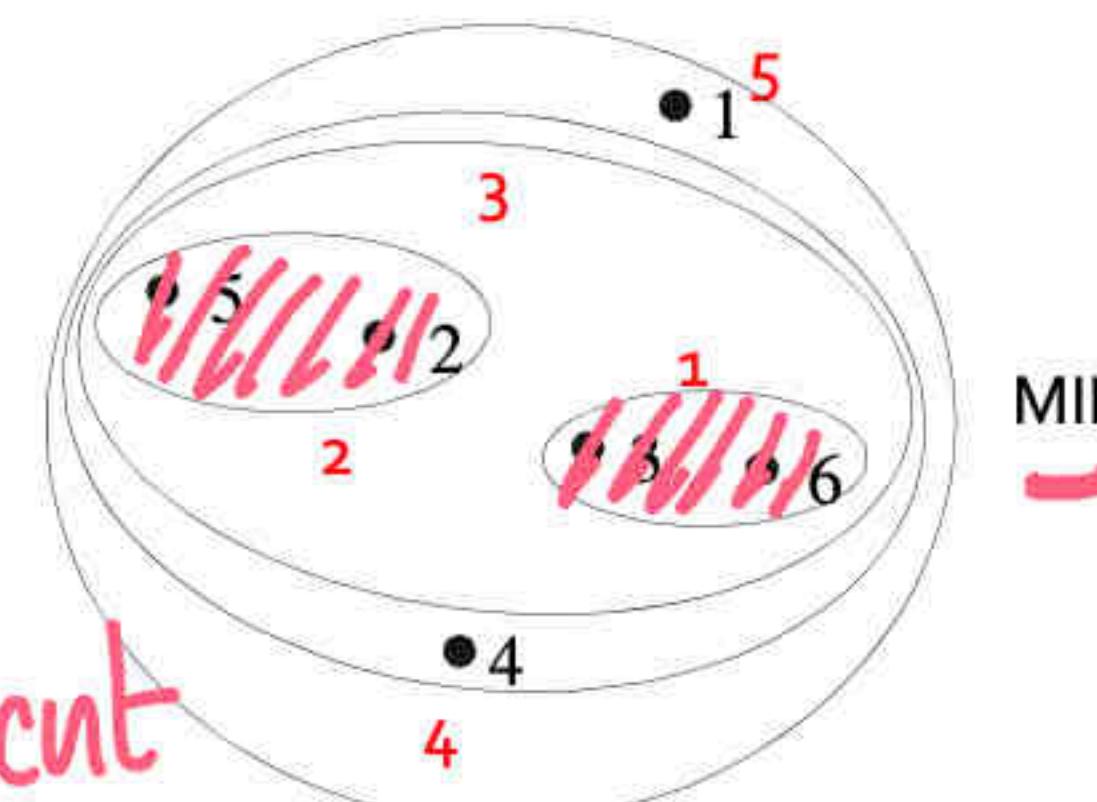
Nested Clusters

Dendrogram

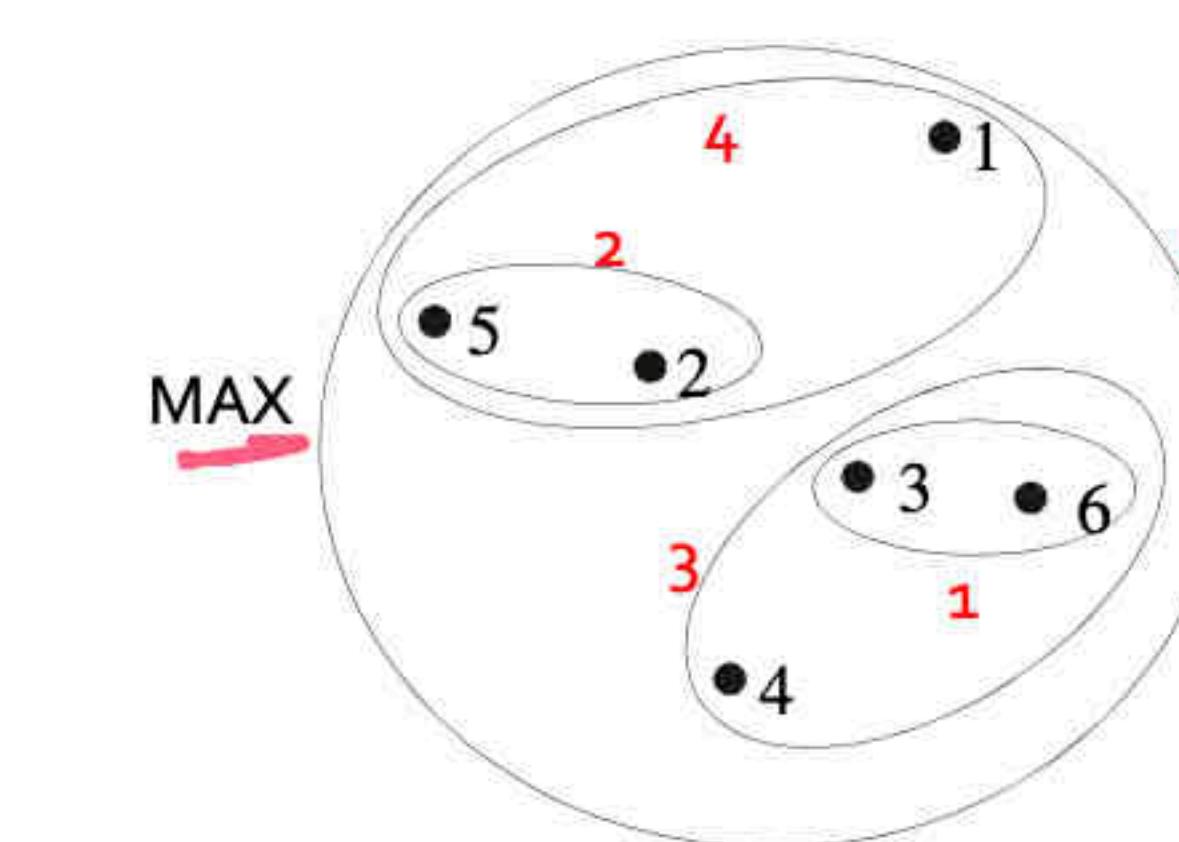
	1	2	3	4	5	6
1	0	.24	.22	.37	.34	.23
2	.24	0	.15	.20	.14	.25
3	.22	.15	0	.15	.28	.11
4	.37	.20	.15	0	.29	.22
5	.34	.14	.28	.29	0	.39
6	.23	.25	.11	.22	.39	0



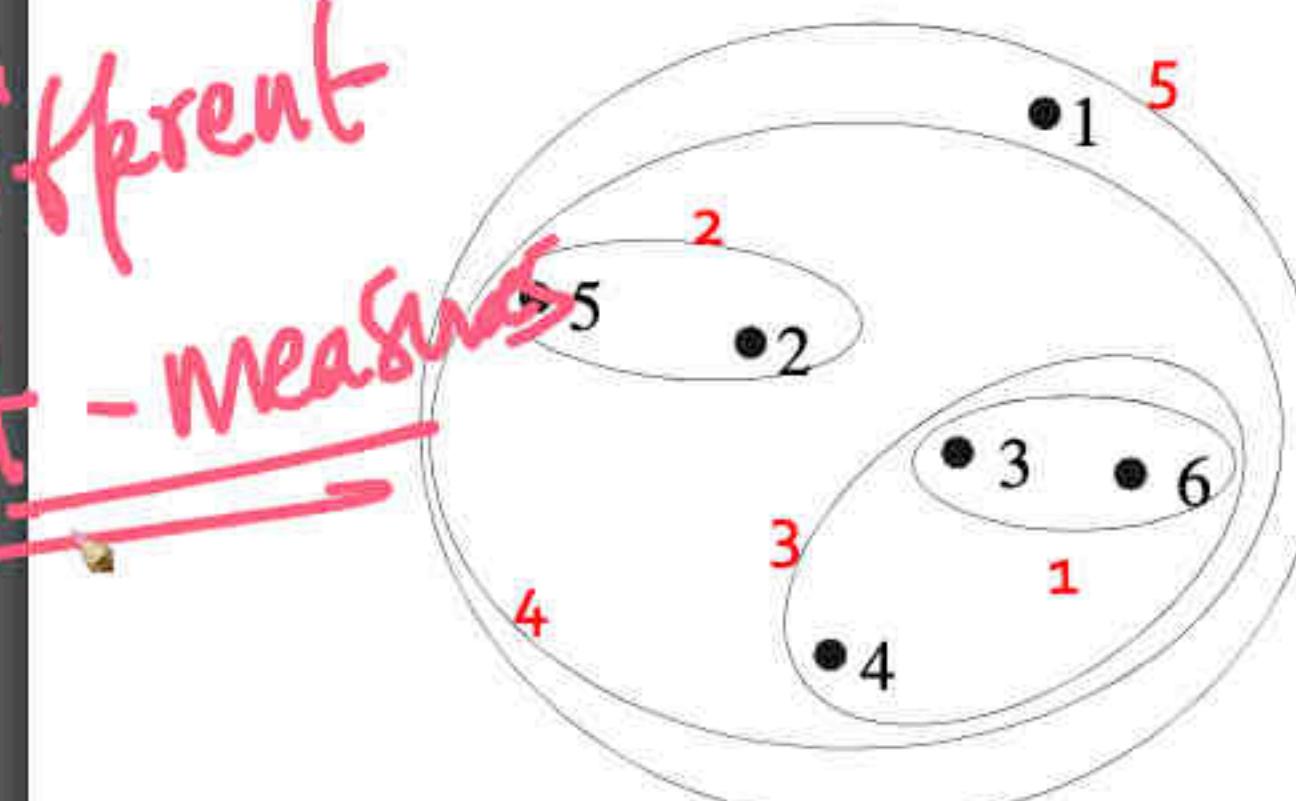
clustering  
results  
can be different  
for different  
dist-measures



MIN

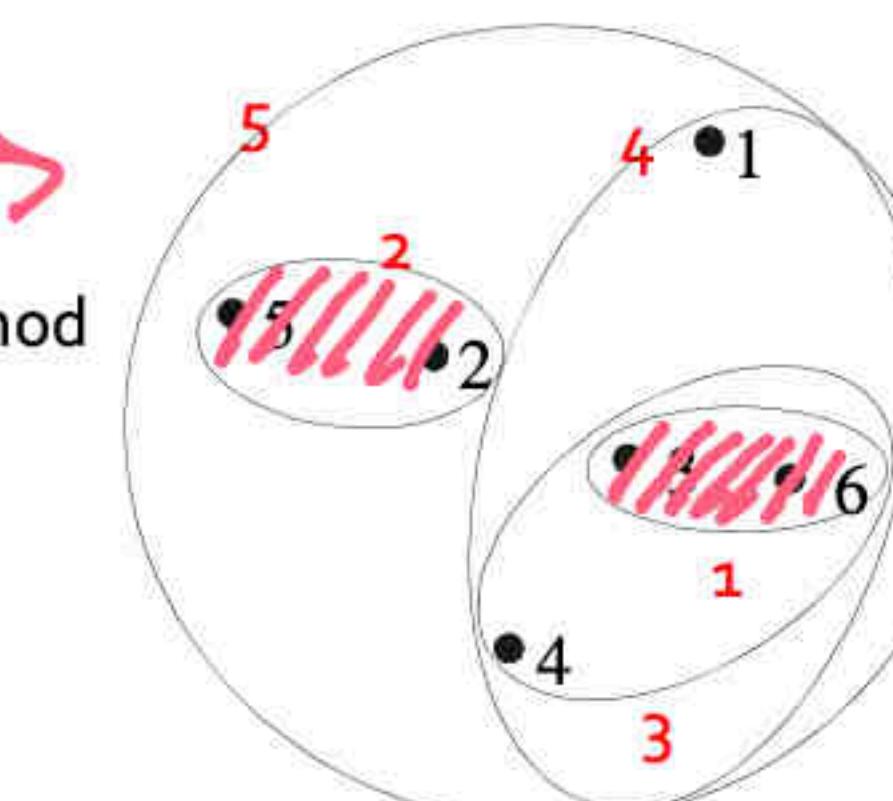


MAX



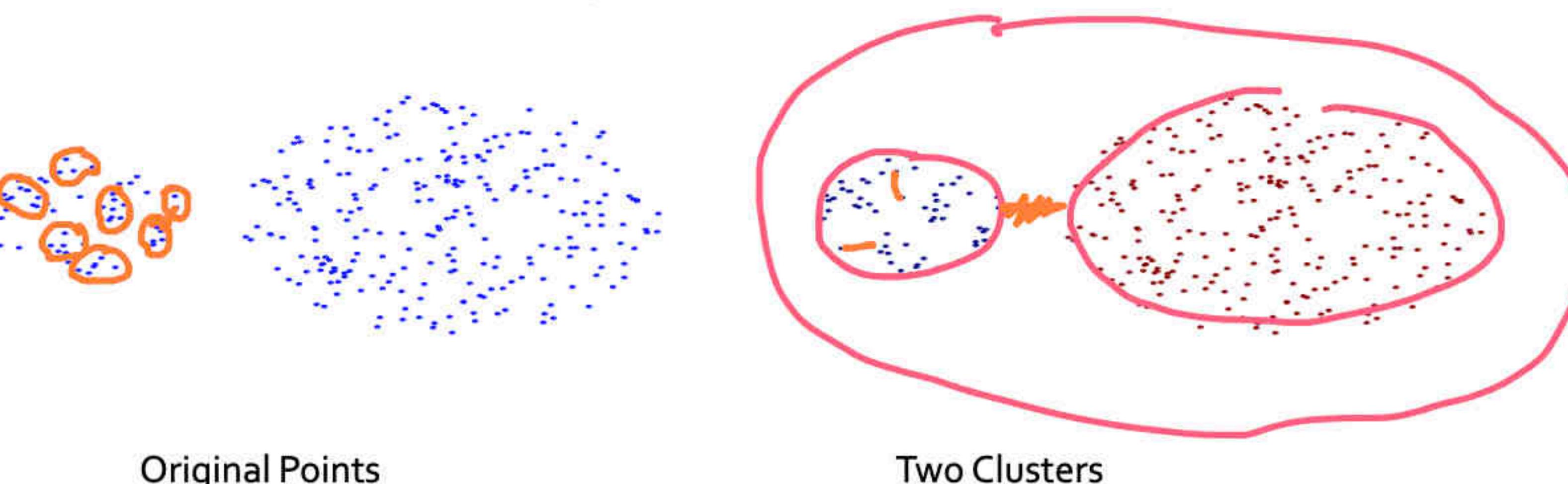
Group Average

Ward's Method  
→



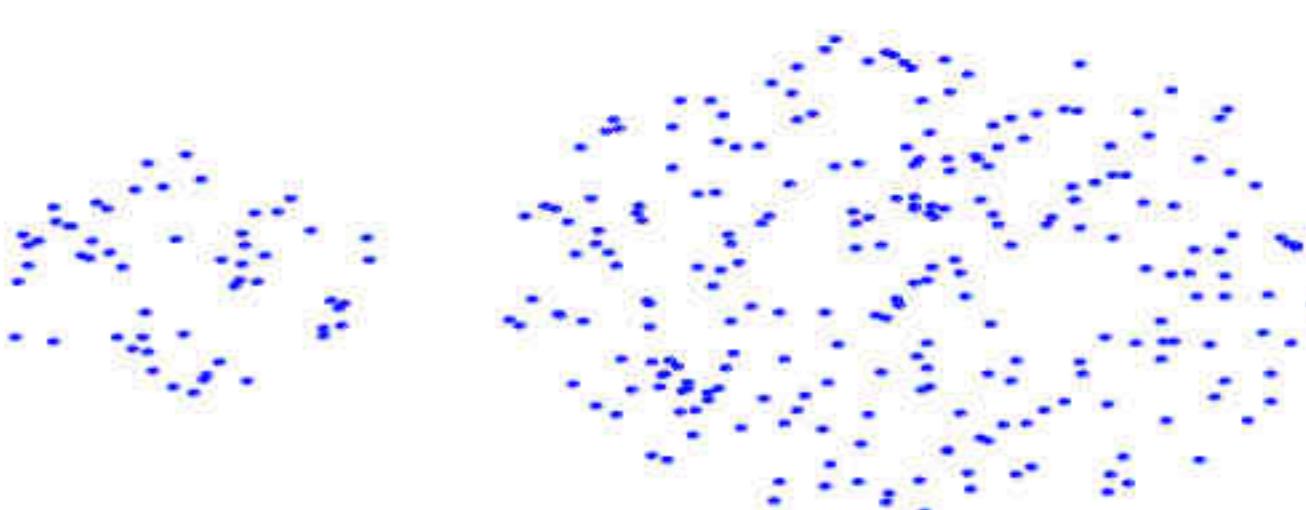
# Hierarchical Clustering: Comparison

# Strength of MIN

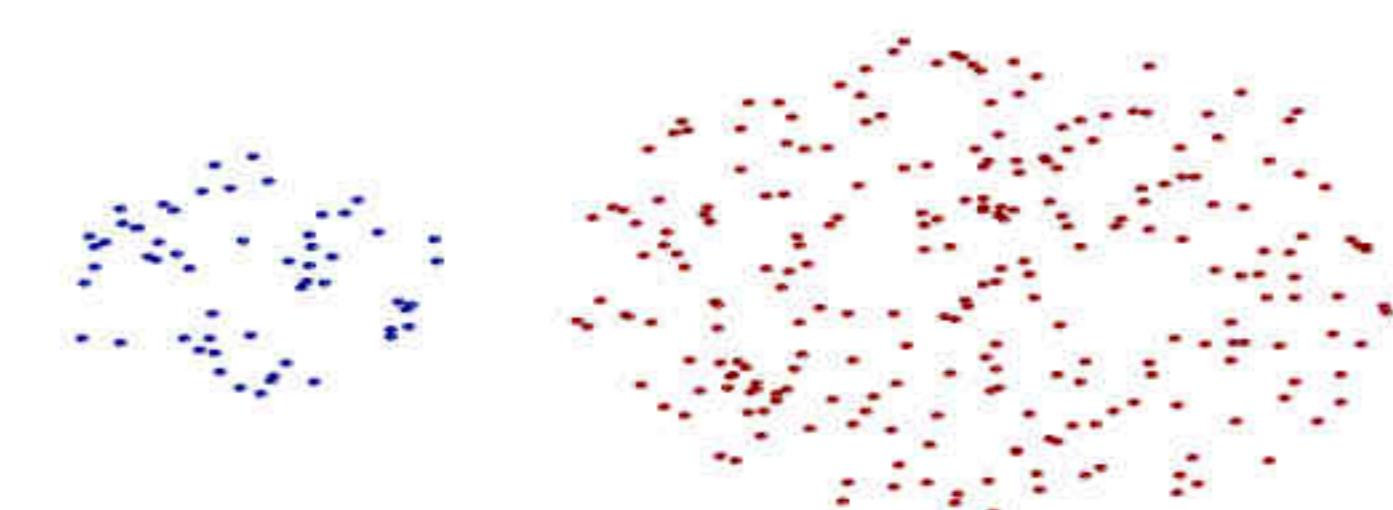


- Can handle non-elliptical shapes

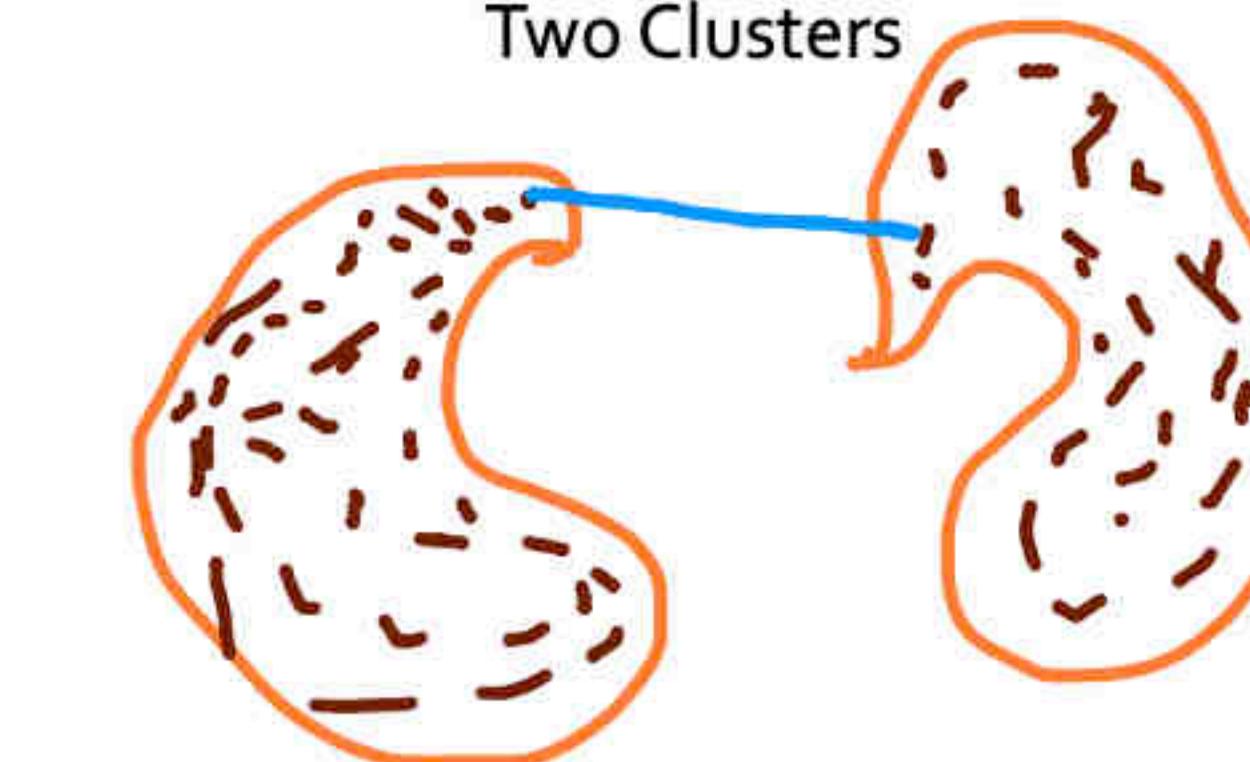
# Strength of MIN



Original Points



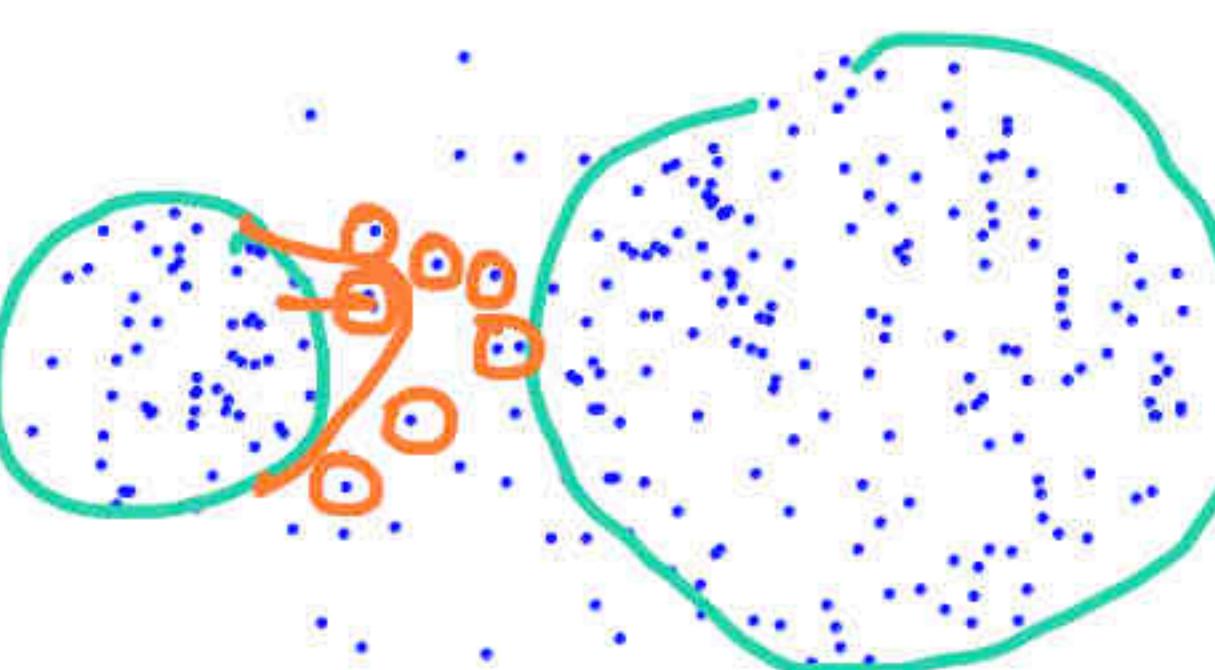
Two Clusters



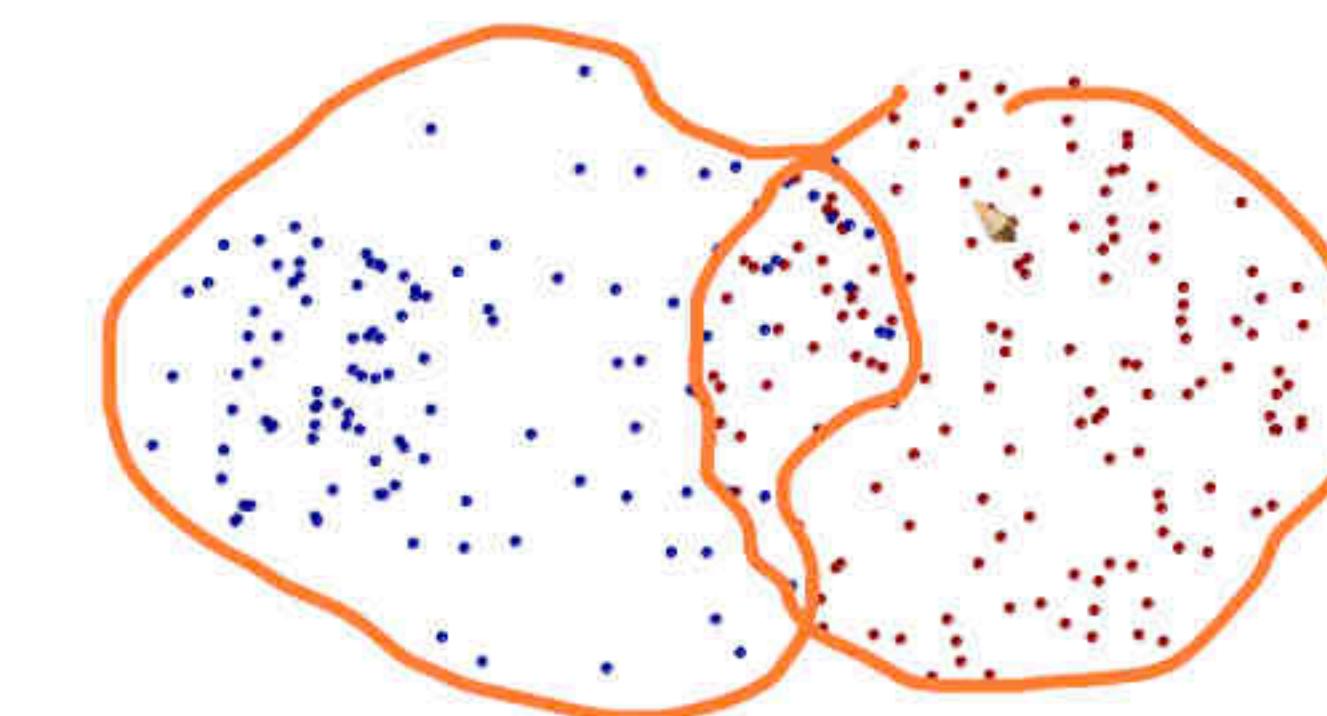
- Can handle non-elliptical shapes

non-globular

# Limitations of MIN



Original Points

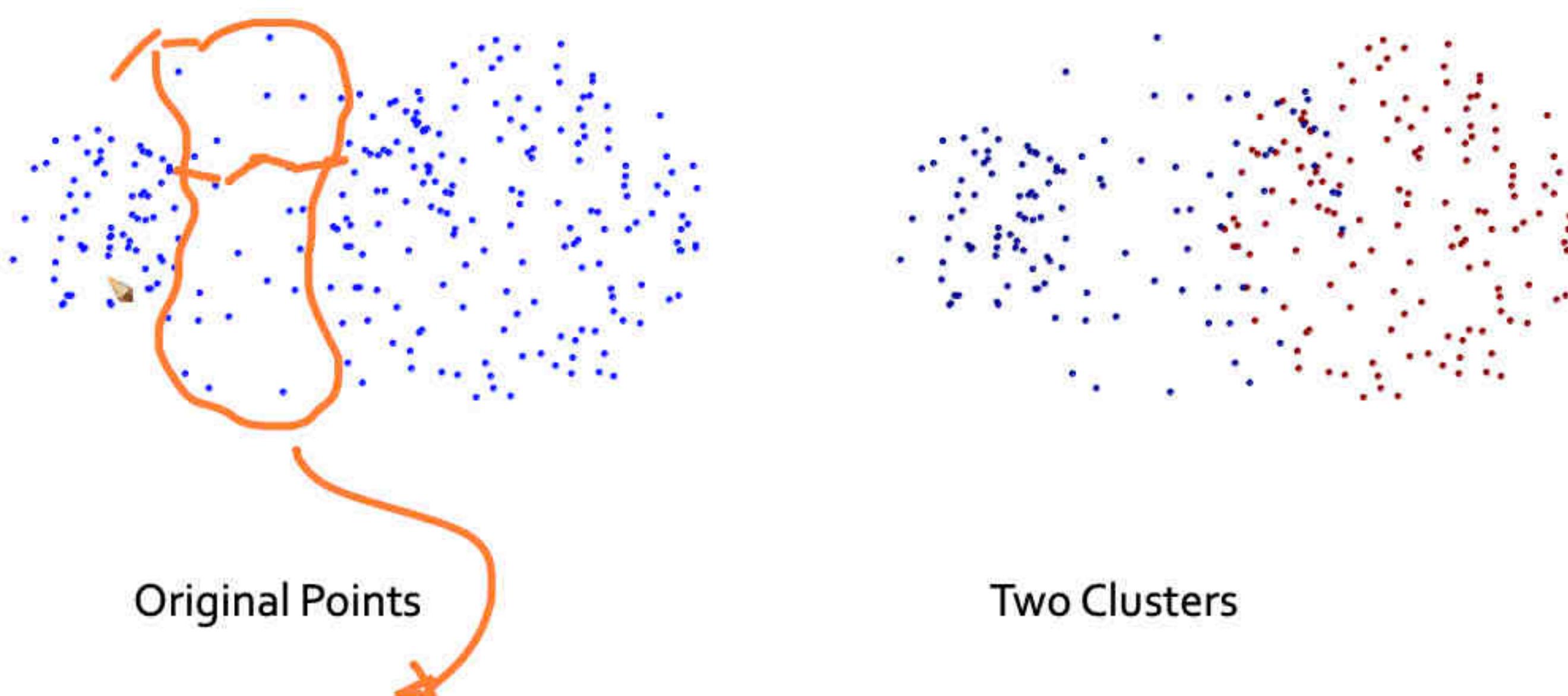


Two Clusters

- Sensitive to noise and outliers

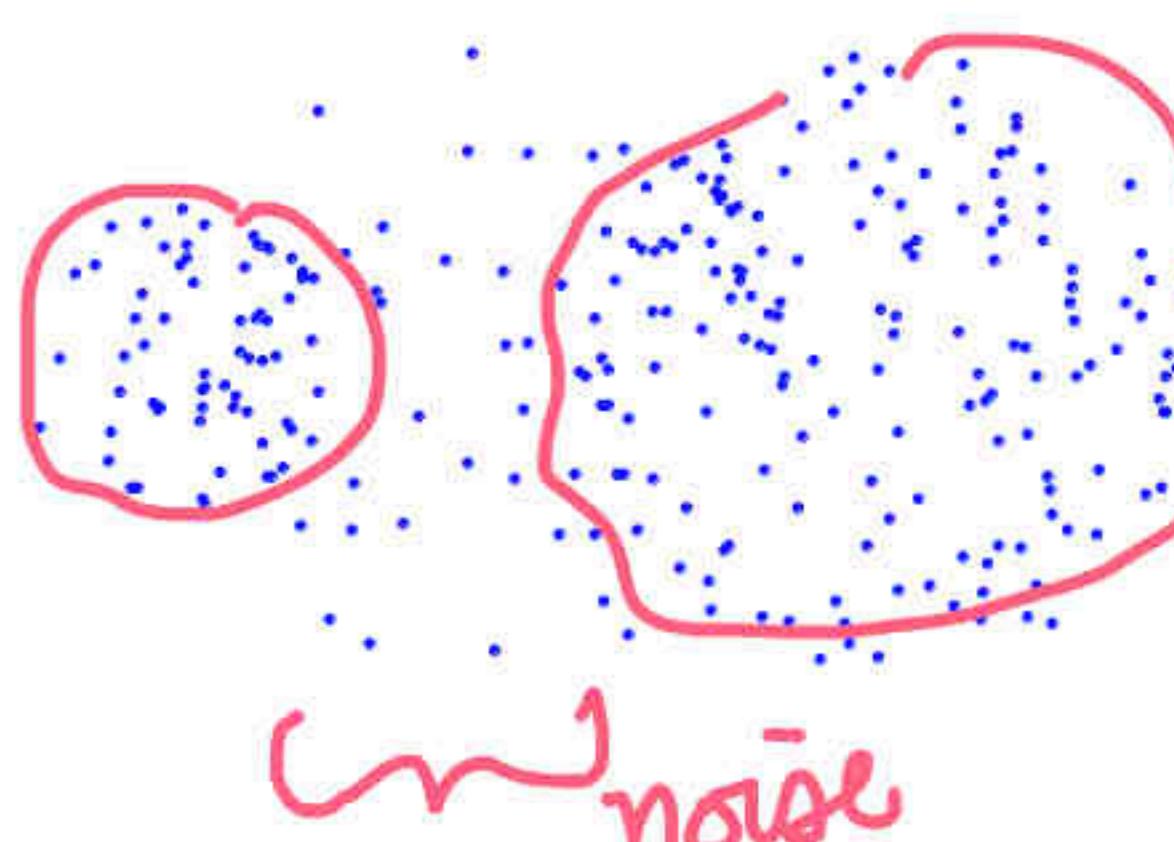


# Limitations of MIN → Single-linkage

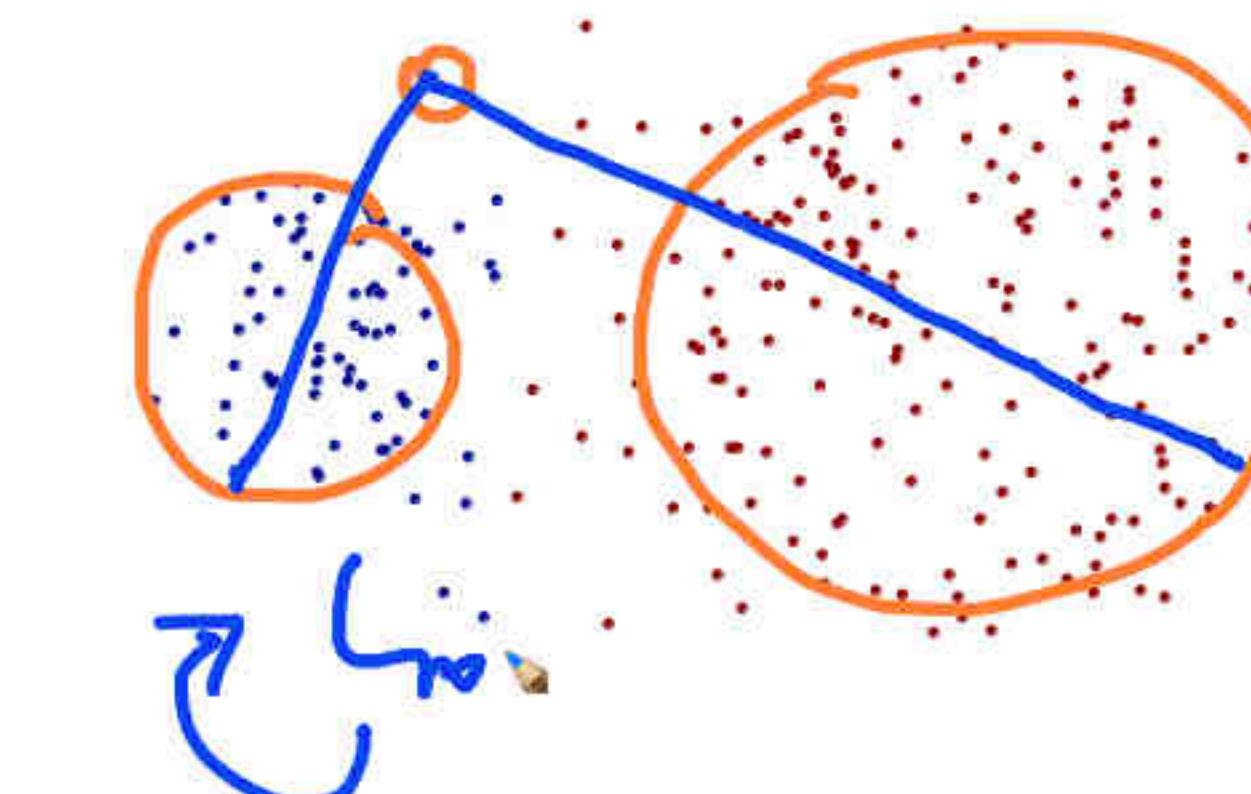


- Sensitive to noise and outliers

# Strength of MAX : complete linkage



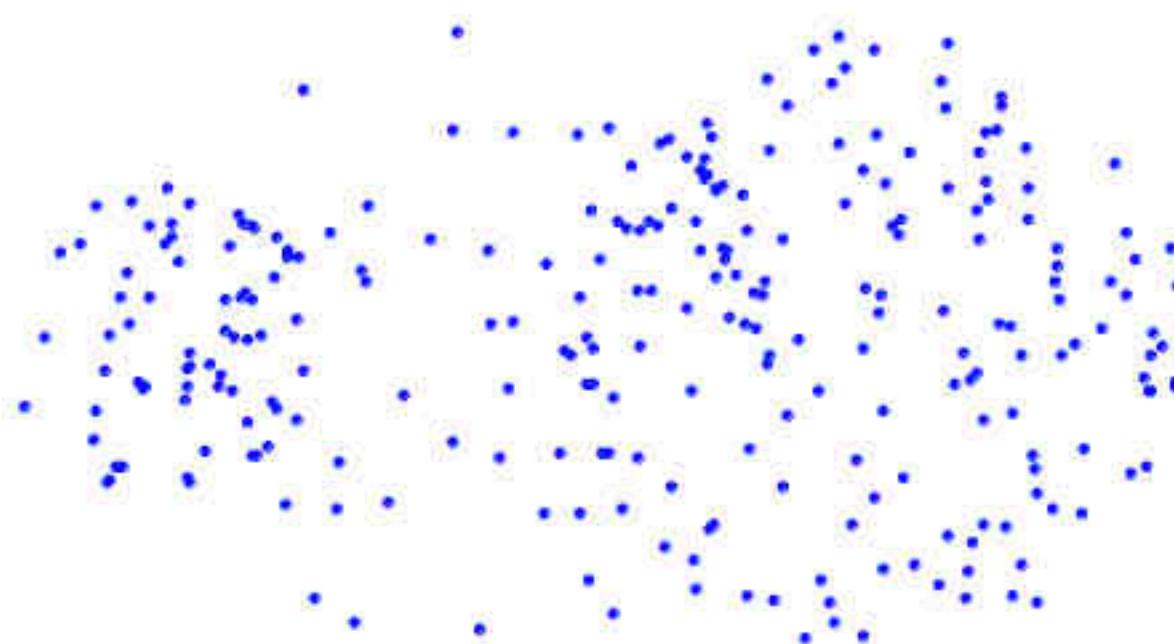
Original Points



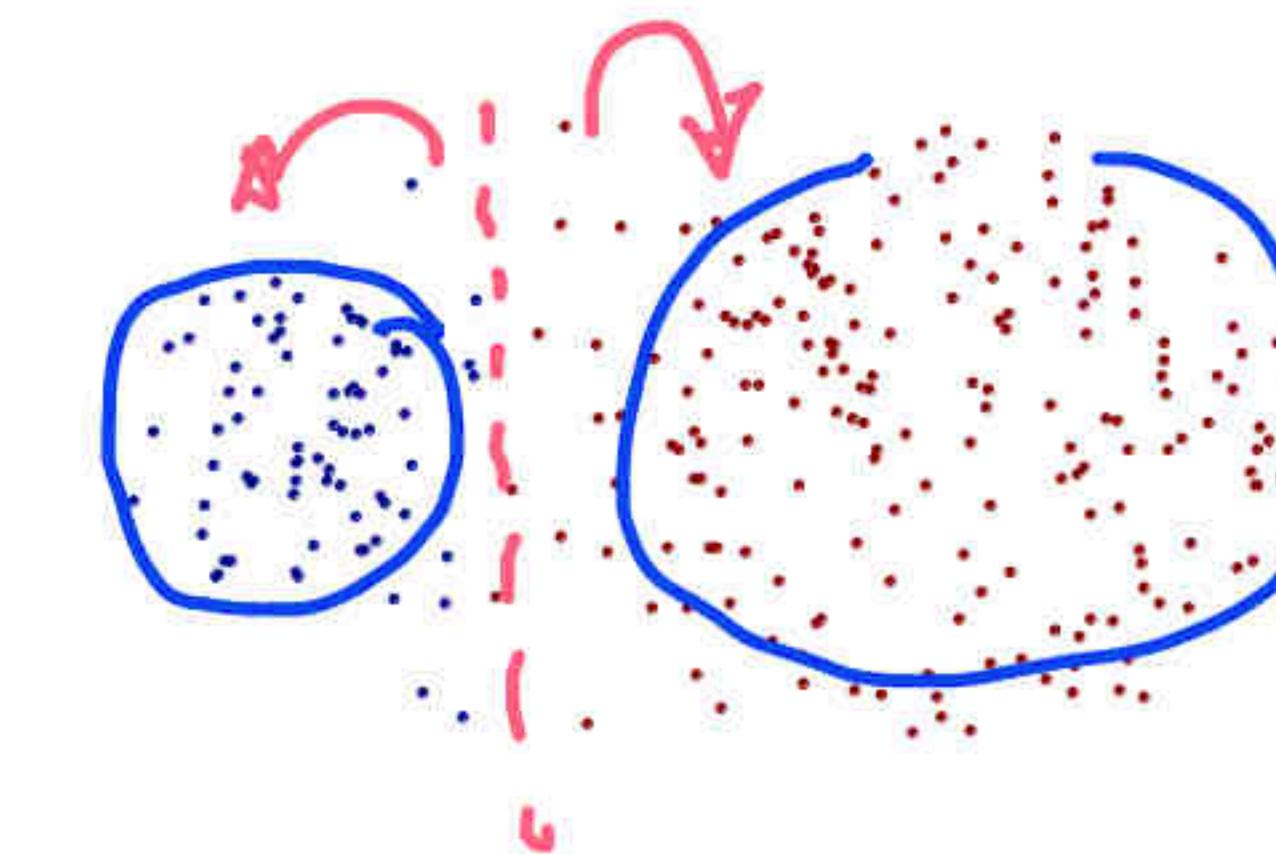
Two Clusters

- Less susceptible to noise and outliers

# Strength of MAX



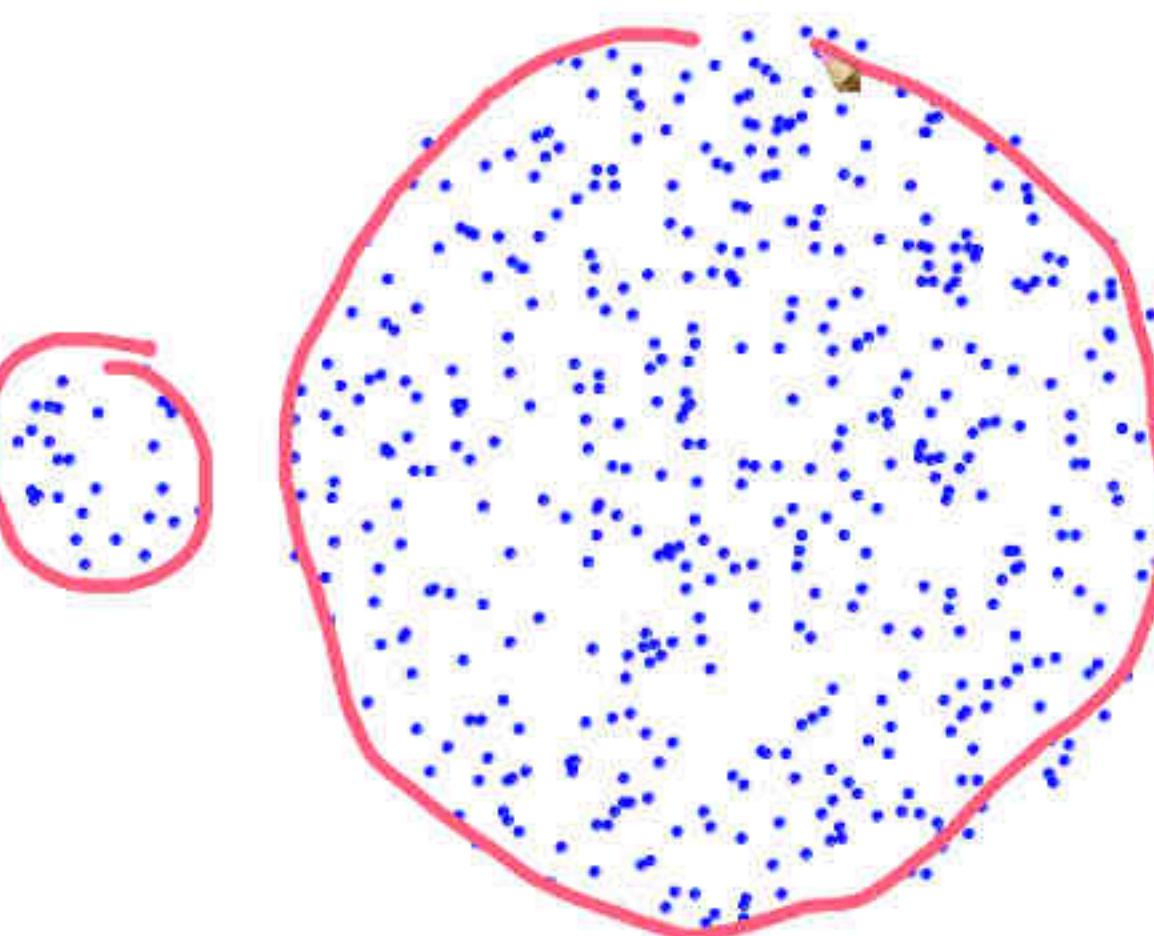
Original Points



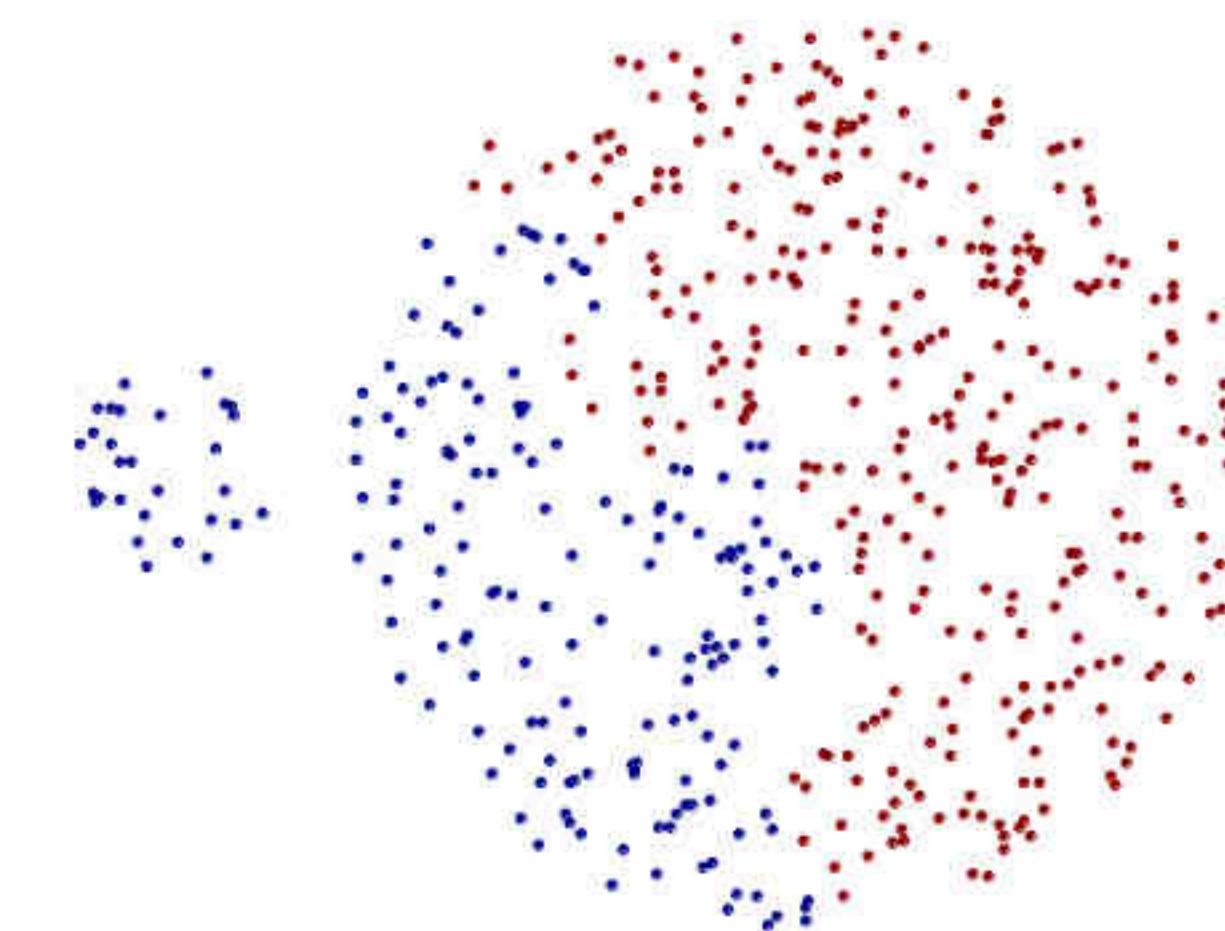
Two Clusters

- Less ~~susceptible~~ to noise and outliers

# Limitations of MAX



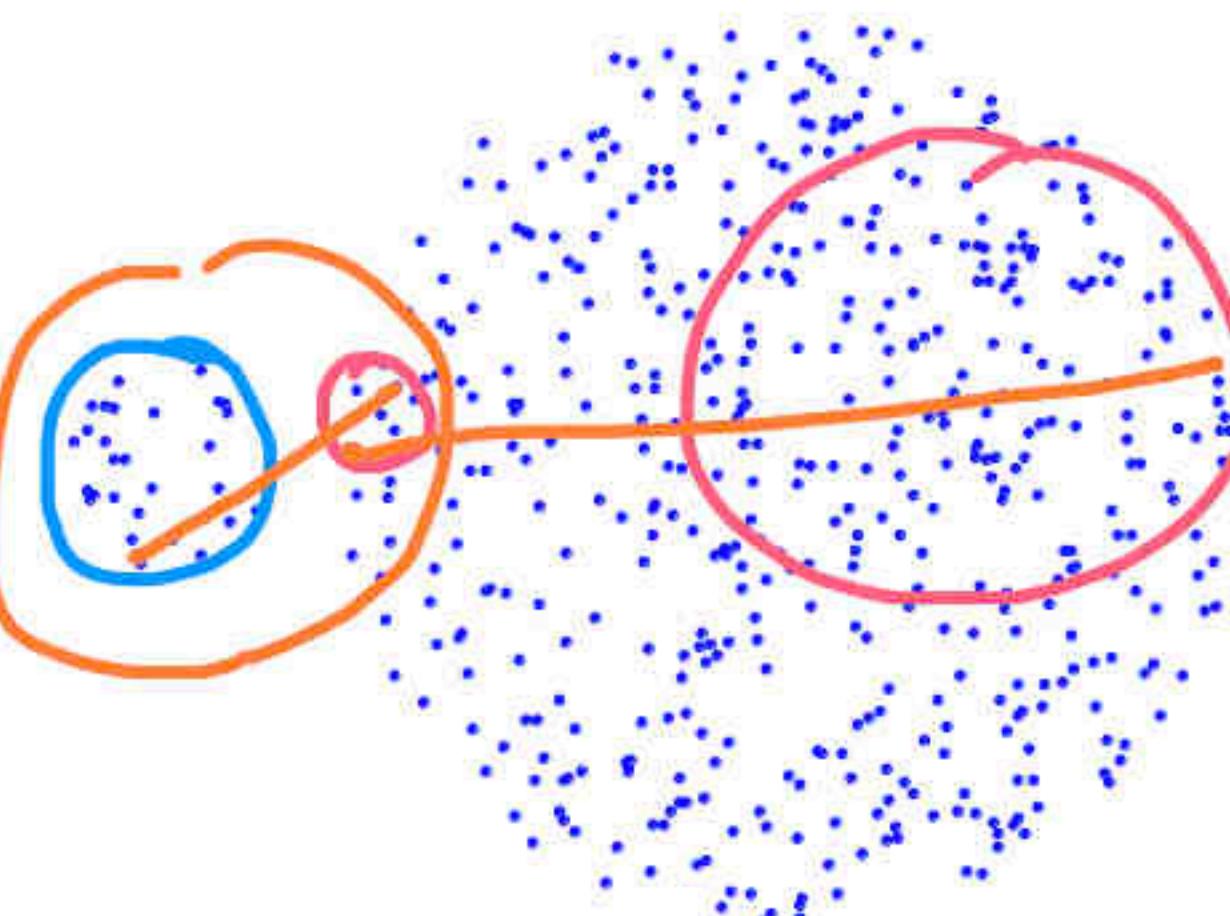
Original Points



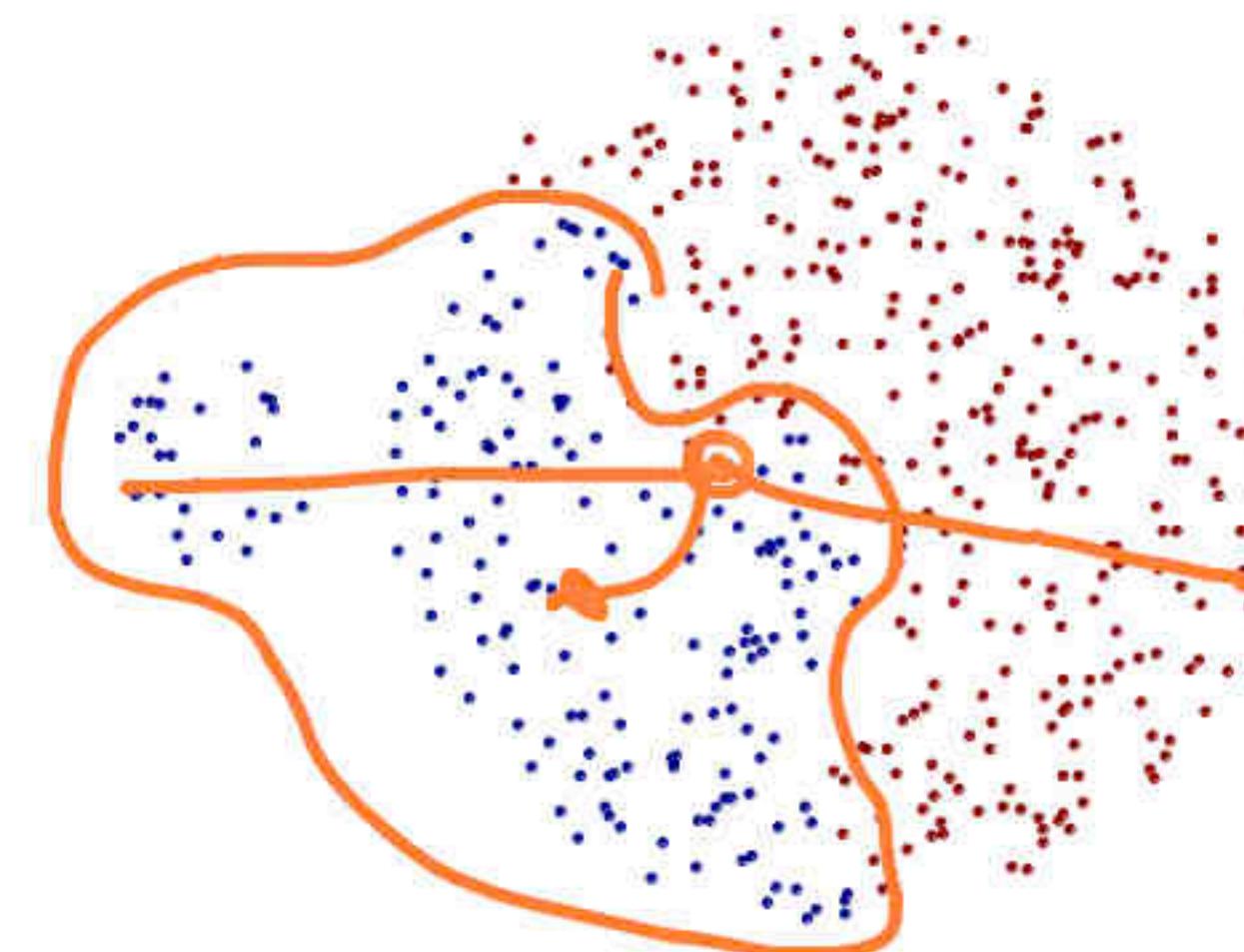
Two Clusters

- Tends to break large clusters
- Biased towards globular clusters

# Limitations of MAX



Original Points



Two Clusters

- Tends to break large clusters
- Biased towards globular clusters

}  $\rightarrow$  k-Means

# Cluster Similarity: Group Average

- Proximity of two clusters is the average of pairwise proximity between points in the two clusters.

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

- Need to use average connectivity for scalability since total proximity favors large clusters

	1	2	3	4	5	6
1	0	.24	.22	.37	.34	.23
2	.24	0	.15	.20	.14	.25
3	.22	.15	0	.15	.28	.11
4	.37	.20	.15	0	.29	.22
5	.34	.14	.28	.29	0	.39
6	.23	.25	.11	.22	.39	0

Group-avg or ward algo

# Hierarchical Clustering: Group Average

- Compromise between Single and Complete Link

Min      Max



- Strengths
  - Less susceptible to noise and outliers
- Limitations
  - Biased towards globular clusters



# Average

- Compromise between Single and Complete Link
- Strengths
  - Less susceptible to noise and outliers
- Limitations
  - Biased towards globular clusters



Microsoft PowerPoint - Cluster X | DBSCAN - Wikipedia X | sklearn.cluster.AgglomerativeClustering.html

scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html

Go

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

**sklearn.cluster.AgglomerativeClustering**

Examples using `sklearn.cluster.AgglomerativeClustering`

# sklearn.cluster.AgglomerativeClustering

```
class sklearn.cluster.AgglomerativeClustering(n_clusters=2, *, affinity='euclidean',  
memory=None, connectivity=None, compute_full_tree='auto', linkage='ward', distance_threshold=None,  
compute_distances=False)
```

[source]

Agglomerative Clustering.

Recursively merges pair of clusters of sample data; uses linkage distance.

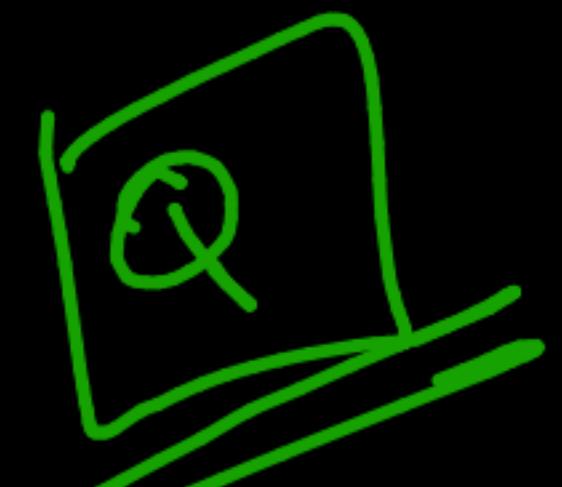
Read more in the [User Guide](#).

**Parameters:** `n_clusters : int or None, default=2`  
The number of clusters to find. It must be `None` if `distance_threshold` is not `None`.

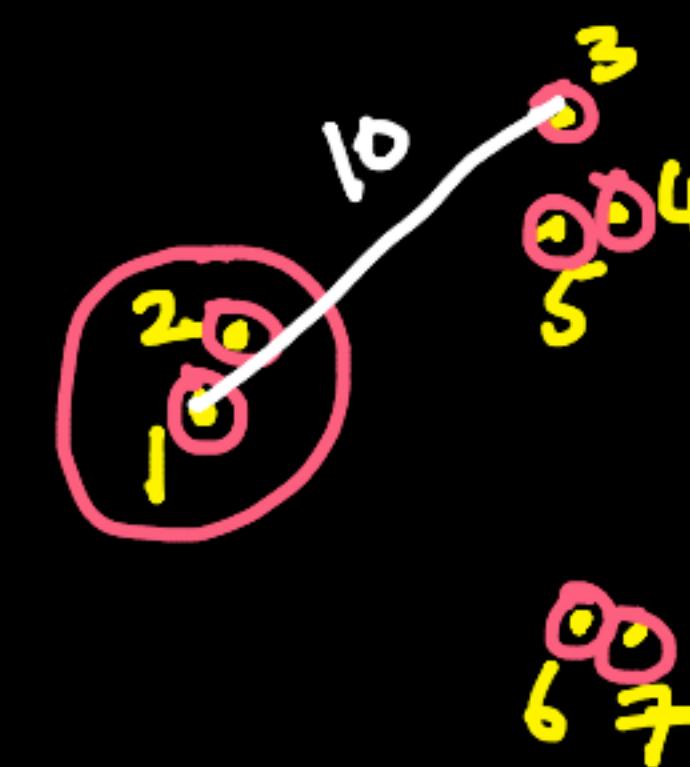
`affinity : str or callable, default='euclidean'`  
Metric used to compute the linkage. Can be "euclidean", "l1", "l2", "manhattan", "cosine", or "precomputed". If linkage is "ward", only "euclidean" is accepted. If "precomputed", a distance matrix (instead of a similarity matrix) is needed as input for the fit method.

Toggle Menu

memory : str or object with the `__reduce_ex__` interface, default=None



MAX / complete linkage



common &  
dangerous  
misconception

P: Min / Max / avg

✓ merge ~~2 closest~~  
clusters

l,2      3,4,5,6,7

l,2	0	10					
3	10	0					
4		0					
5			0				
6				0			
7					0		

$n$  is large  $\rightarrow$

Agglomerative  
clustering

Space:  $O(n \times n)$ : Proximity-matrix

~~Expensive~~

Time:  $\sim O(n^3)$   
worst-case

optimized: -  $O(n^2 \lg n)$

$n$  [ - compute  $P$   
repeat { merge; update  $P$  }  
until one cluster ]  $n^2 \times$

$n \rightarrow n-1 \rightarrow n-2 \rightarrow \dots 1$

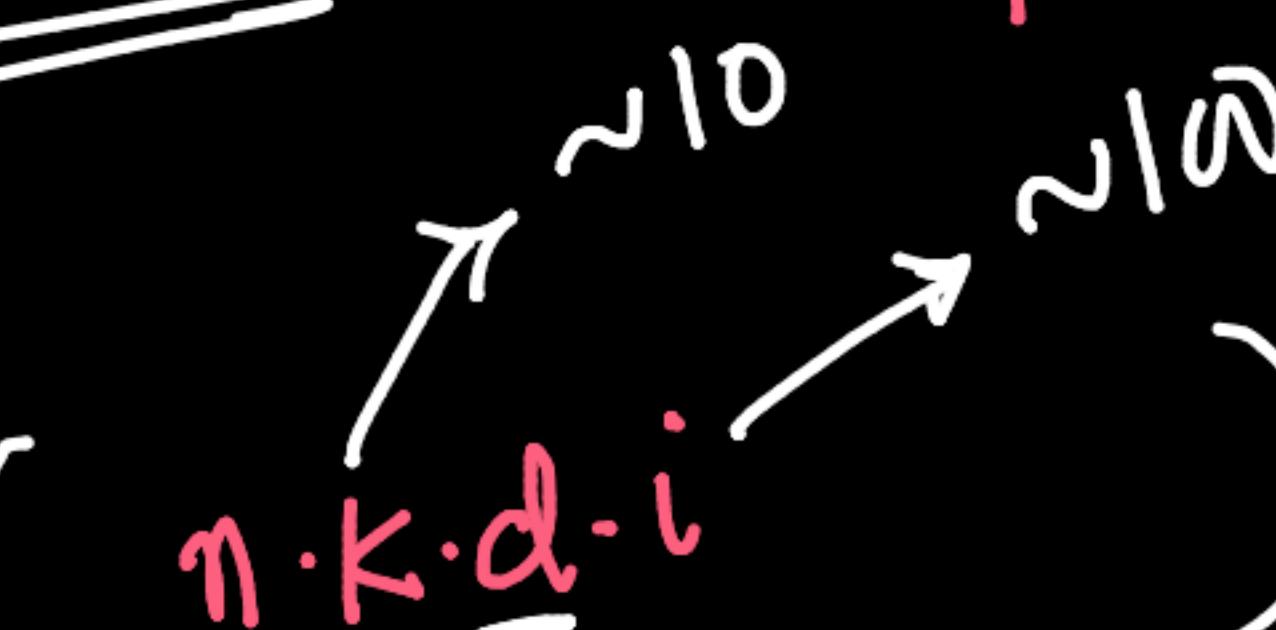
K-means:

$$\underline{\underline{O(nd)}}$$

→ Computationally Cheaper

Time:

$$\underline{\underline{O(n \cdot k \cdot d \cdot i)}}$$



$\stackrel{\sim}{=} 100$   
i iterations



$\stackrel{\sim}{=} 10$   
 $d\text{-dim } j \stackrel{\sim}{=} K$

$k \ll n$

$k \ll d$

$i \ll n$

- {}  $\stackrel{\sim}{=} O(k)$   
random centroids
- {} for  $i$ -iterations
  - {}  $\stackrel{\sim}{=} n \cdot k \cdot d$
  - {} → assign each pt to a cluster
- {} → update centroids
- {}  $\stackrel{\sim}{=} O(nd)$

No Math obj that is being minimized in  
Agg-clustering

K-means:  $\min \underline{\underline{\text{within cluster ss-dist}}}$

MIN & MAX  $\rightarrow$  resolved using Ward's method

Microsoft PowerPoint - Cluster X | DBSCAN - Wikipedia X | sklearn.cluster.AgglomerativeClustering.html +

scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html

Go



Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.cluster.AgglomerativeClustering](#)

Examples using

[sklearn.cluster.AgglomerativeClustering](#)

## sklearn.cluster.AgglomerativeClustering

```
class sklearn.cluster.AgglomerativeClustering(n_clusters=2, *, affinity='euclidean',
memory=None, connectivity=None, compute_full_tree='auto', linkage='ward', distance_threshold=None,
compute_distances=False) [source]
```

Agglomerative Clustering.

Recursively merges pair of clusters of sample data; uses linkage distance.

Read more in the [User Guide](#).

**Parameters:** **n\_clusters : int or None, default=2**

The number of clusters to find. It must be `None` if `distance_threshold` is not `None`.

**affinity : str or callable, default='euclidean'**

Metric used to compute the linkage. Can be "euclidean", "l1", "l2", "manhattan", "cosine", or "precomputed". If linkage is "ward", only "euclidean" is accepted. If "precomputed", a distance matrix (instead of a similarity matrix) is needed as input for the fit method.



Microsoft PowerPoint - Cluster X | DBSCAN - Wikipedia X | sklearn.cluster.AgglomerativeClustering.html +

scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html

Read more in the [User Guide](#).

**Parameters:**

- n\_clusters : int or None, default=2**  
The number of clusters to find. It must be `None` if `distance_threshold` is not `None`.
- affinity : str or callable, default='euclidean'**  
Metric used to compute the linkage. Can be "euclidean", "l1", "l2", "manhattan", "cosine", or "precomputed". If linkage is "ward", only "euclidean" is accepted. If "precomputed", a distance matrix (instead of a similarity matrix) is needed as input for the fit method.
- memory : str or object with the joblib.Memory interface, default=None**  
Used to cache the output of the computation of the tree. By default, no caching is done. If a string is given, it is the path to the caching directory.
- connectivity : array-like or callable, default=None**  
Connectivity matrix. Defines for each sample the neighboring samples following a given structure of the data. This can be a connectivity matrix itself or a callable that transforms the data into a connectivity matrix, such as derived from `kneighbors_graph`. Default is `None`, i.e., the hierarchical clustering algorithm is unstructured.
- compute\_full\_tree : 'auto' or bool, default='auto'**

stars. This is useful to decrease computation time if the number of clusters is not small compared to the number of

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.cluster.AgglomerativeClustering

Examples using sklearn.cluster.AgglomerativeClustering

0(n<sup>2</sup>) Space

Toggle Menu

Microsoft PowerPoint - Cluster X | DBSCAN - Wikipedia X | sklearn.cluster.AgglomerativeClustering.html +

scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html

CONNECTIVITY : array-like or callable, default=None

Connectivity matrix. Defines for each sample the neighboring samples following a given structure of the data. This can be a connectivity matrix itself or a callable that transforms the data into a connectivity matrix, such as derived from `kneighbors_graph`. Default is `None`, i.e., the hierarchical clustering algorithm is unstructured.

**compute\_full\_tree : 'auto' or bool, default='auto'** = 10

Stop early the construction of the tree at `n_clusters`. This is useful to decrease computation time if the number of clusters is not small compared to the number of samples. This option is useful only when specifying a connectivity matrix. Note also that when varying the number of clusters and using caching, it may be advantageous to compute the full tree. It must be `True` if `distance_threshold` is not `None`. By default `compute_full_tree` is "auto", which is equivalent to `True` when `distance_threshold` is not `None` or that `n_clusters` is inferior to the maximum between 100 or  $0.02 * n\_samples$ . Otherwise, "auto" is equivalent to `False`.

**linkage : {'ward', 'complete', 'average', 'single'}, default='ward'**

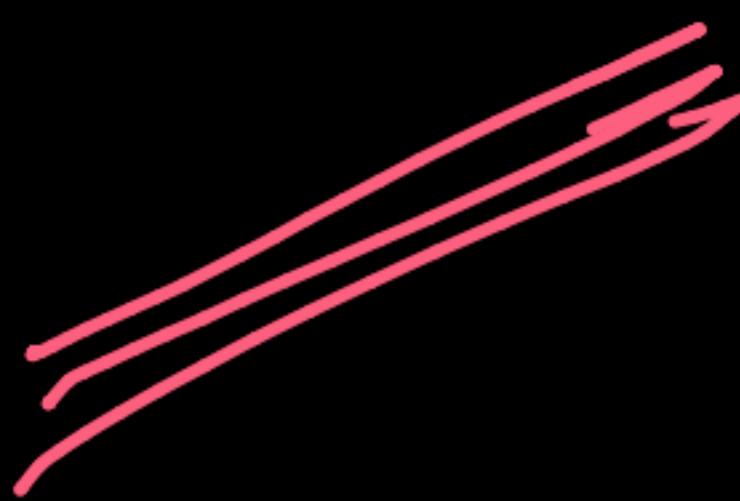
Which linkage criterion to use. The linkage criterion determines which distance to use between sets of observation. The algorithm will merge the pairs of cluster that minimize this criterion.

- 'ward' minimizes the variance of the clusters being merged.

of each observation of the two sets.

Toggle Menu





K-means: Centroid - based

Agg.-clustering: hierarchical system

DB-SCAN → density - based

databases,  
data mining

WIKIPEDIA  
The Free Encyclopedia[Main page](#)  
[Contents](#)  
[Current events](#)  
[Random article](#)  
[About Wikipedia](#)  
[Contact us](#)  
[Donate](#)  
[Contribute](#)[Help](#)  
[Learn to edit](#)  
[Community portal](#)  
[Recent changes](#)  
[Upload file](#)[Tools](#)  
[What links here](#)  
[Related changes](#)  
[Special pages](#)  
[Permanent link](#)  
[Page information](#)  
[Cite this page](#)  
[Wikidata item](#)Article [Talk](#)Read [Edit](#) [View history](#)

Search Wikipedia



# DBSCAN

From Wikipedia, the free encyclopedia

**Density-based spatial clustering of applications with noise (DBSCAN)** is a [data clustering](#) algorithm proposed by [Martin Ester](#), [Hans-Peter Kriegel](#), [Jörg Sander](#) and [Xiaowei Xu](#) in 1996.<sup>[1]</sup> It is a [density-based clustering](#) non-parametric algorithm: given a set of points in some space, it groups together points that are closely packed together (points with many [nearby neighbors](#)), marking as outliers points that lie alone in low-density regions (whose nearest neighbors are too far away). DBSCAN is one of the most common clustering algorithms and also most cited in scientific literature.<sup>[2]</sup>

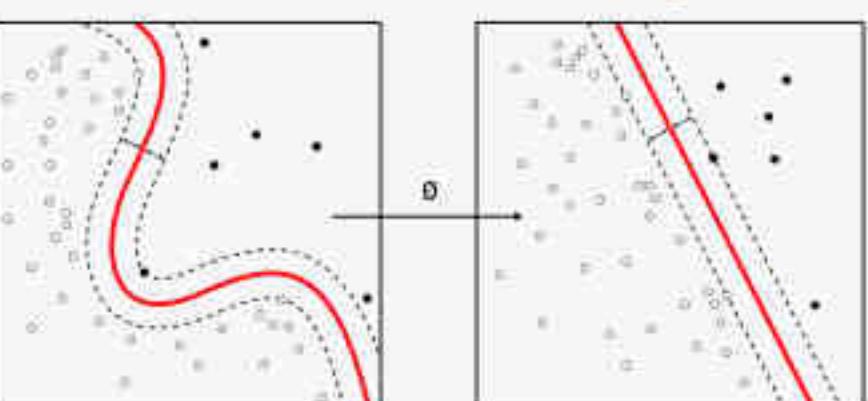
In 2014, the algorithm was awarded the test of time award (an award given to algorithms which have received substantial attention in theory and practice) at the leading data mining conference, ACM [SIGKDD](#).<sup>[3]</sup> As of July 2020, the follow-up paper "DBSCAN Revisited, Revisited: Why and How You Should (Still) Use DBSCAN"<sup>[4]</sup> appears in the list of the 8 most downloaded articles of the prestigious [ACM Transactions on Database Systems \(TODS\)](#) journal.<sup>[5]</sup>

## Contents [hide]

- 1 History
- 2 Preliminary
- 3 Algorithm
  - 3.1 Original query-based algorithm
  - 3.2 Abstract algorithm
- 4 Complexity
- 5 Advantages
- 6 Disadvantages
- 7 Parameter estimation
- 8 Relationship to spectral clustering



Part of a series on  
**Machine learning  
and data mining**



[Problems](#) [show]

[Supervised learning](#) [show]  
(classification · regression)

[Clustering](#) [show]

BIRCH · CURE · Hierarchical ·  $k$ -means ·  
Expectation–maximization (EM) ·  
DBSCAN · OPTICS · Mean shift

[Dimensionality reduction](#) [show]

[Structured prediction](#) [show]

[Anomaly detection](#) [show]

[Artificial neural network](#) [show]

[Reinforcement learning](#) [show]



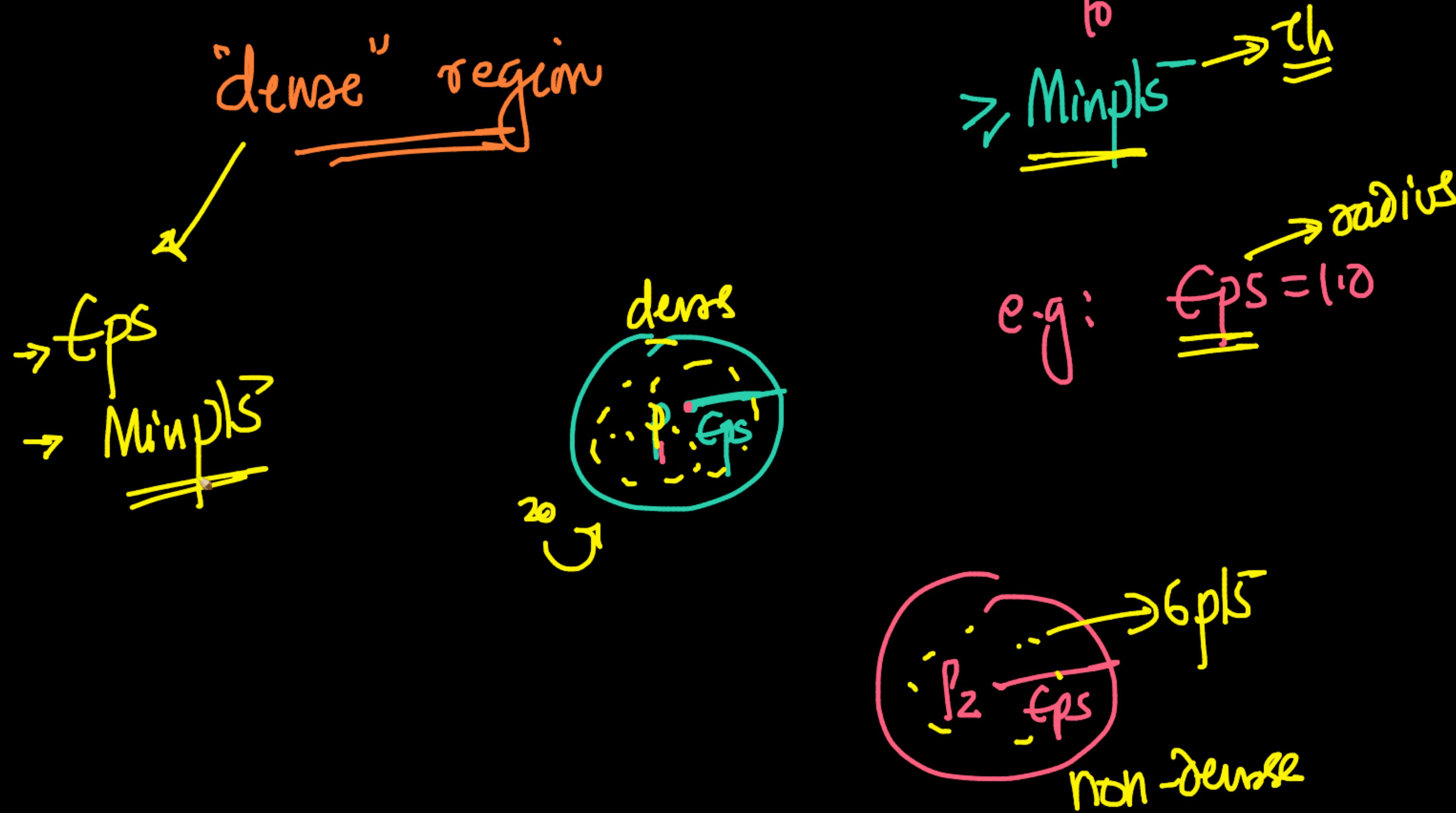
# Terms/concepts (key-ideas)

{ Min pts } → hyper-params  
Eps  
Core-point  
border-point  
noise-point

density @  $\vec{P}$  = # points within the hyperphere  
centered @  $\vec{P}$  with radius =  $\vec{\epsilon}_{PS}$

1.0





✓ { Core-point: ~~if  $P \in S$~~

↳ { if  $P$  has  $\geq M_{\text{pts}}$  in an  $\epsilon_{\text{ps}}$  radius  
around it then  $P$  is a core-pt

$\Downarrow$   $\epsilon_{\text{ps}}, M_{\text{pts}}$

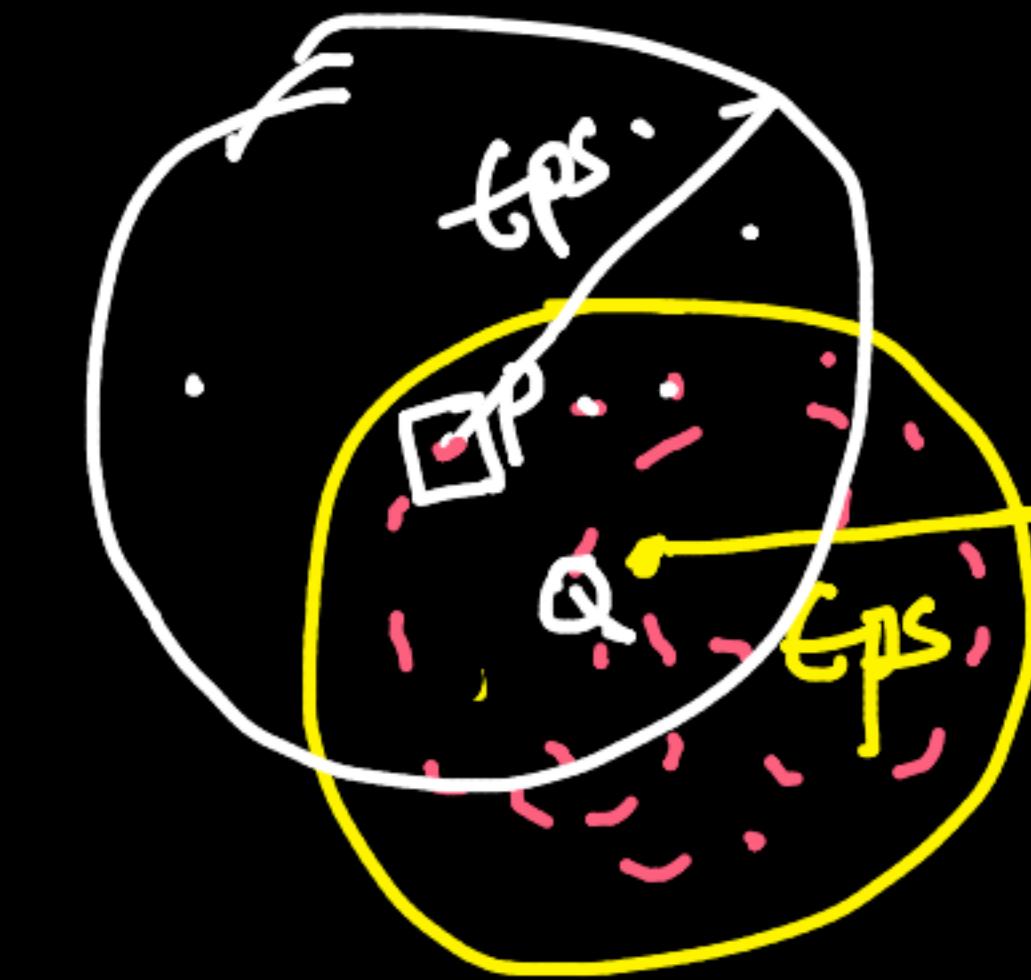
$P$  has a "dense" region around it

border-pt:

→  is not a core-pt

→  $p \in \text{Neighborhood}(q)$  s.t.  $q$  is a core pt

$p \in$  some dense region ||  $\text{dist}(p, q) \leq \epsilon_{\text{ps}}$



not - vole

density @ P = 6  $\neq$  Minpts

density @ Q = 12  $>$  Minpts

↓

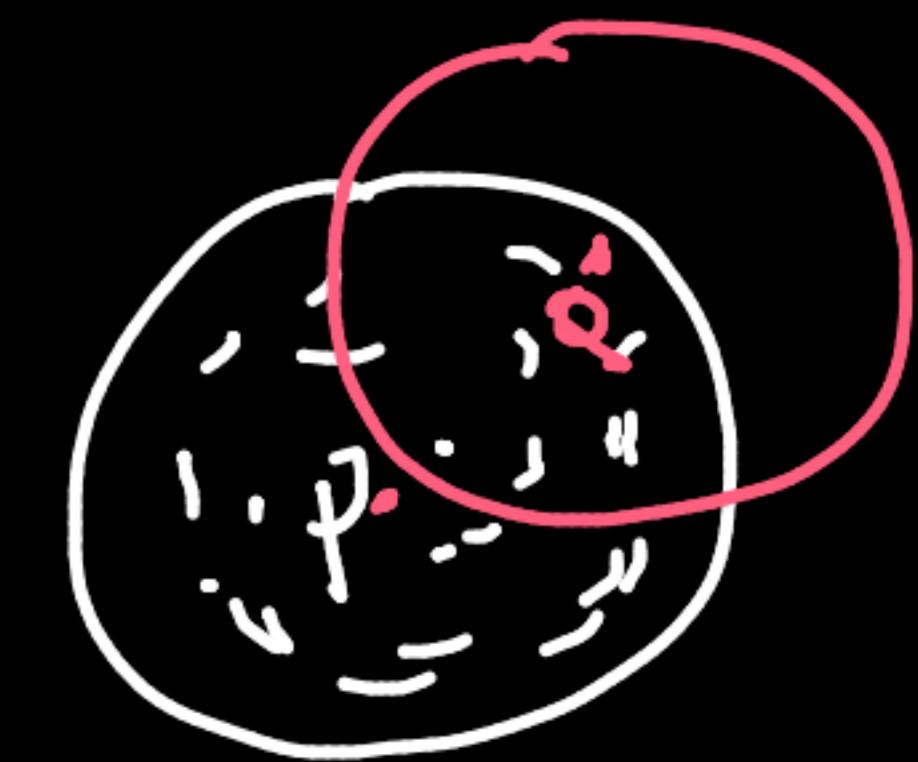
Vole

Minpts = 10

Noise pt:

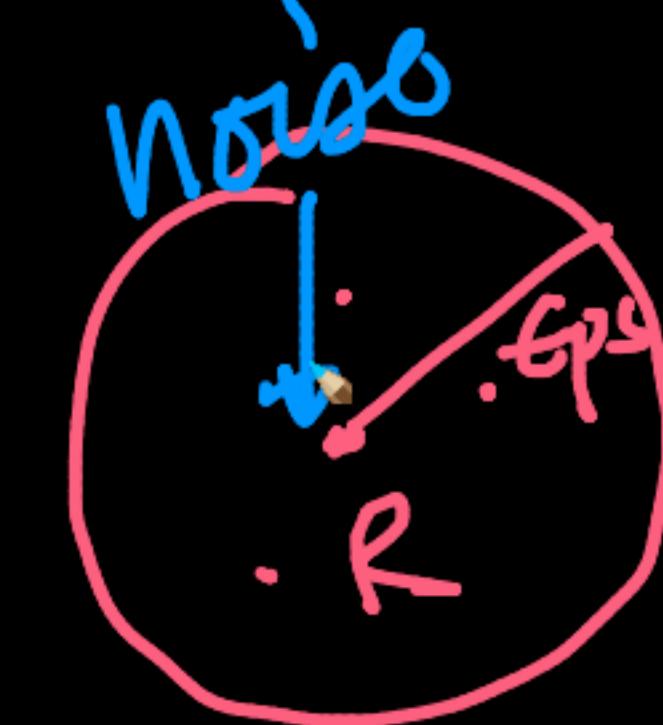
neither core nor

border pt



$p$ : core

$p$ : border



Too few pts  
in  $\epsilon$ -radius

non-core

not  
border

$\rightarrow$   
R not in  
a dense  
region

DB SCAN → Minpts ; Eps → hyper-gism

↓  
core pts ; border pts  
noise pts

Chrome File Edit View History Bookmarks Profiles Tab Window Help

Microsoft PowerPoint - Cluster DBSCAN - Wikipedia sklearn.cluster.Agglomerative 73 / 102 - 125% + | Microsoft PowerPoint - ClusteringAnalysis.pptx cs.wmich.edu/alfuqaha/summer14/cs6530/lectures/ClusteringAnalysis.pdf Update

# DBSCAN: CORE, BORDER and NOISE Points

Original Points

Eps = 10, MinPts = 4

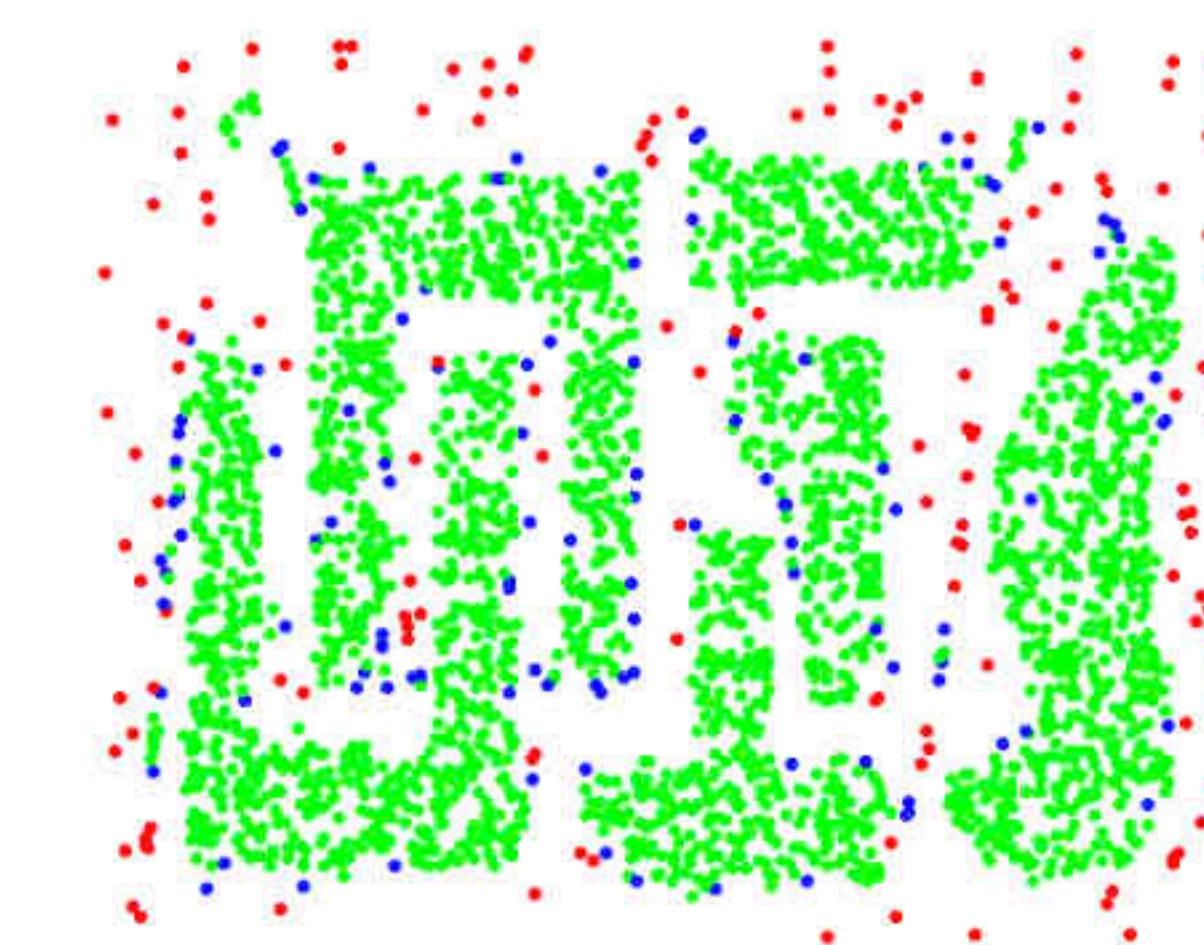
Point types: core, border and noise

46 / 46

# DBSCAN: CORE, BORDER and NOISE Points



Original Points



Point types: core, border and noise

Eps = 10, MinPts = 4

~~Eps = 10, MinPts = 4~~

density-edge/connexum



edge-pq is a density edge

if P & Q are core-points (and)

↳ dense regions

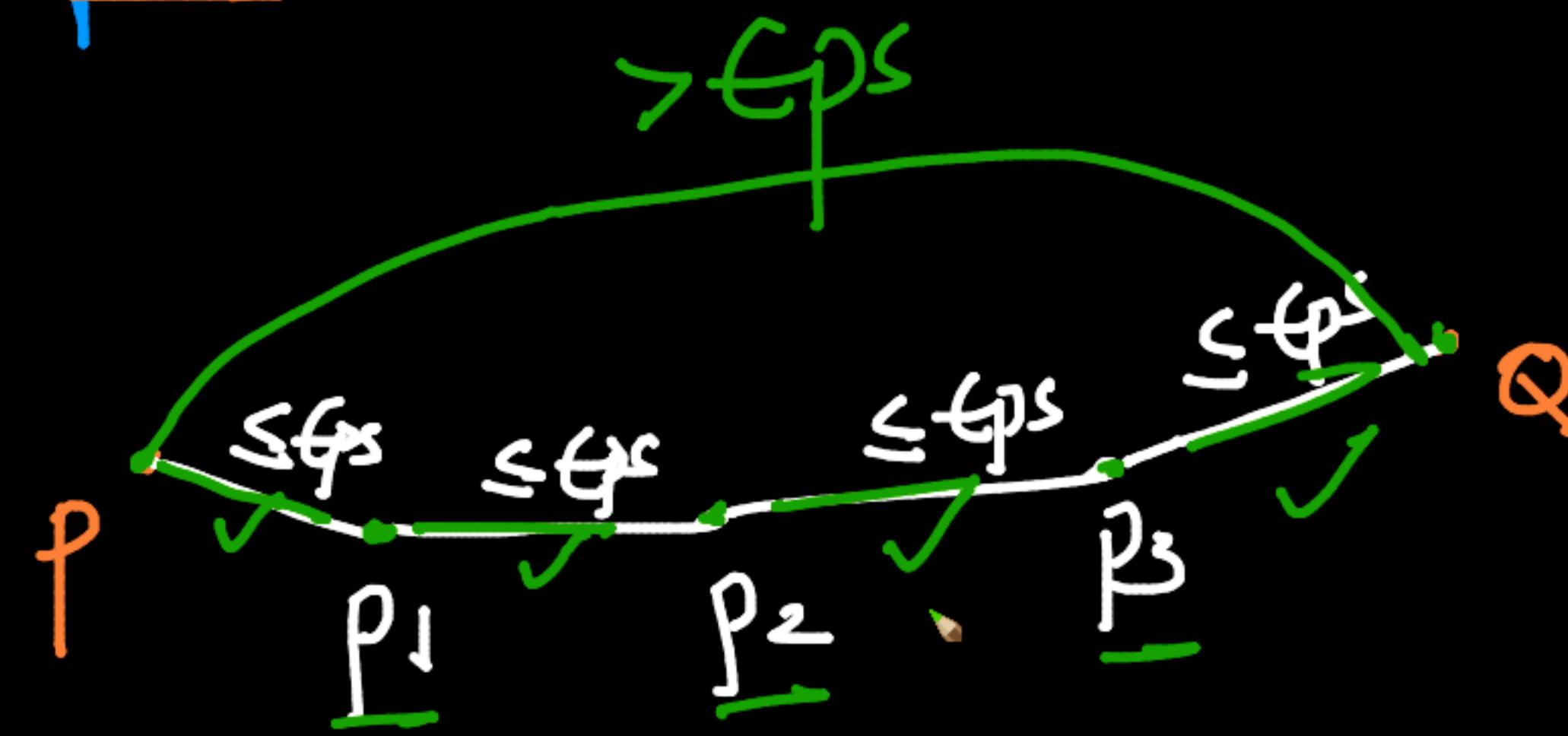


✓ { each-pt is a vertex  
line  $P_1 - P_2$  : edge

$\text{dist}(P, Q) \leq \epsilon_{ps}$

vertex/node  
Graphs

density connected pls



{ P & Q: core pls  
P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>: core pls

P & Q are density connected if there  
exist  $\rightarrow$  density edges connecting P & Q

→ dense regions

{

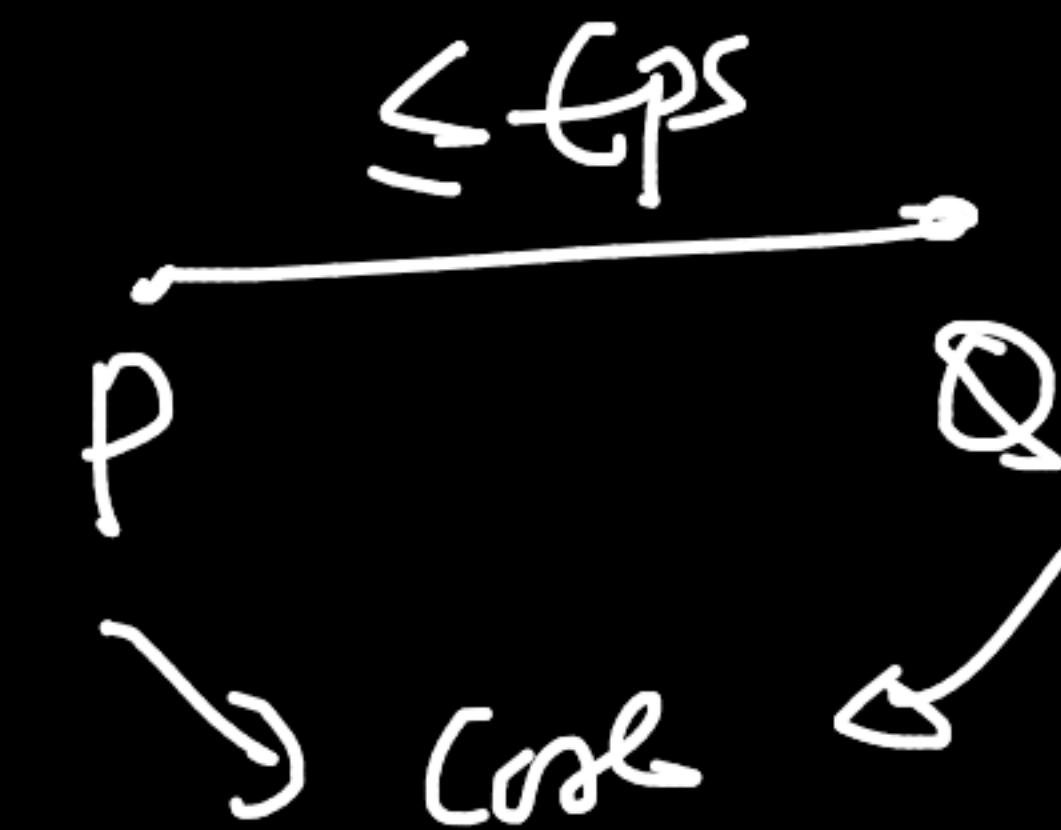
- Minpts,  $\epsilon$  pts
- #pts in  $\epsilon$  pts



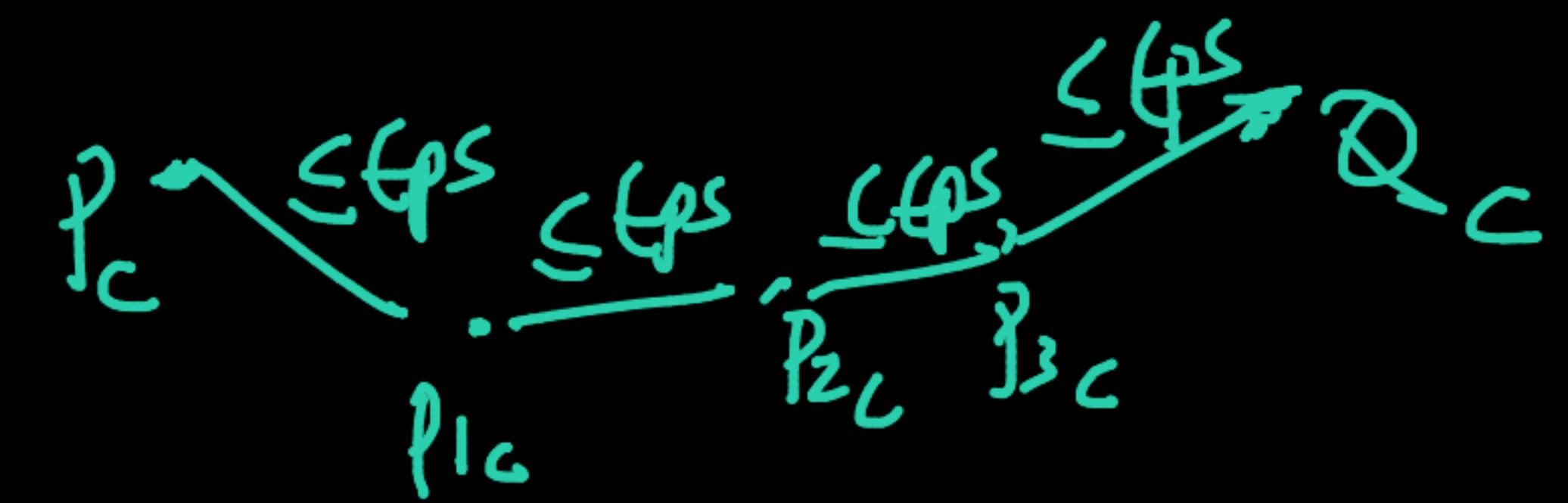
{

- Core pt: dense regions
- border pt: - non core but pts in a dense border
- noise pt: - non-core + non-border

density edge



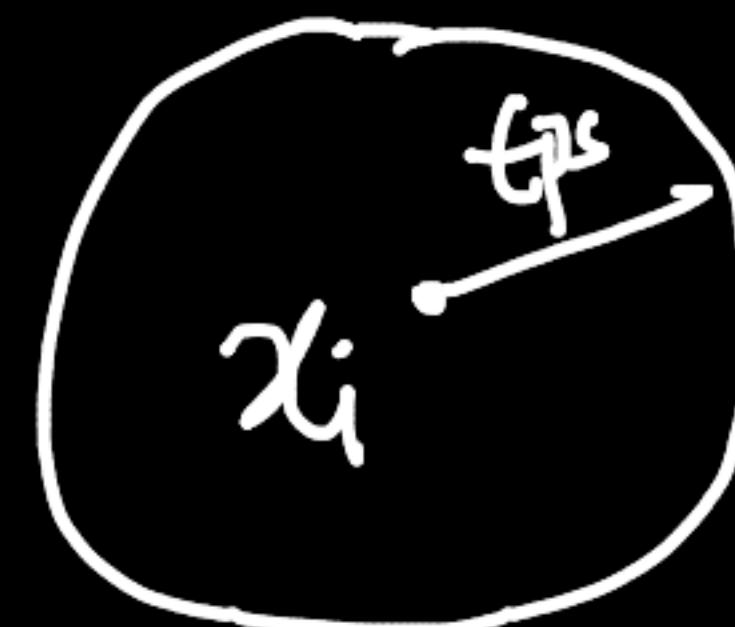
density connected



(Minpts, eps) DB-SCAN:

①  $\forall x_i \in D$

~~$n \lg n$~~   $\rightarrow$



label concept: borderpt; noise pt  
DS: R\*-trees, k-d-trees ...

RangeQuery( $D, x_i, \text{eps}$ )  $\lg(n)$   
#pts / density

noise-removal (density)  
remove all noise pts from your datasets  
Sparse regions

2

$O(n)$

3

$n \cdot \lg n$

For each core pt

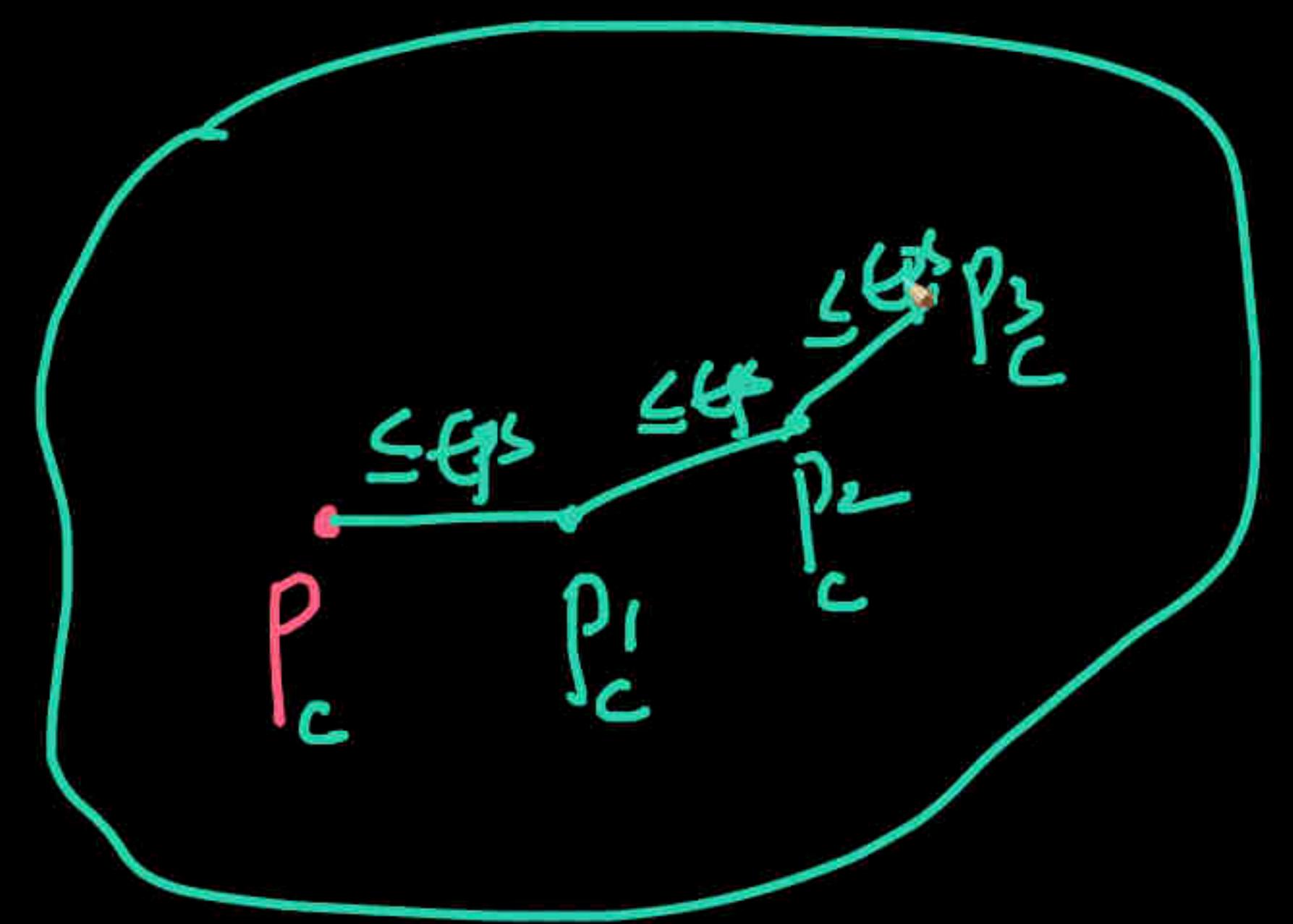
$m$   
P

$P_C$

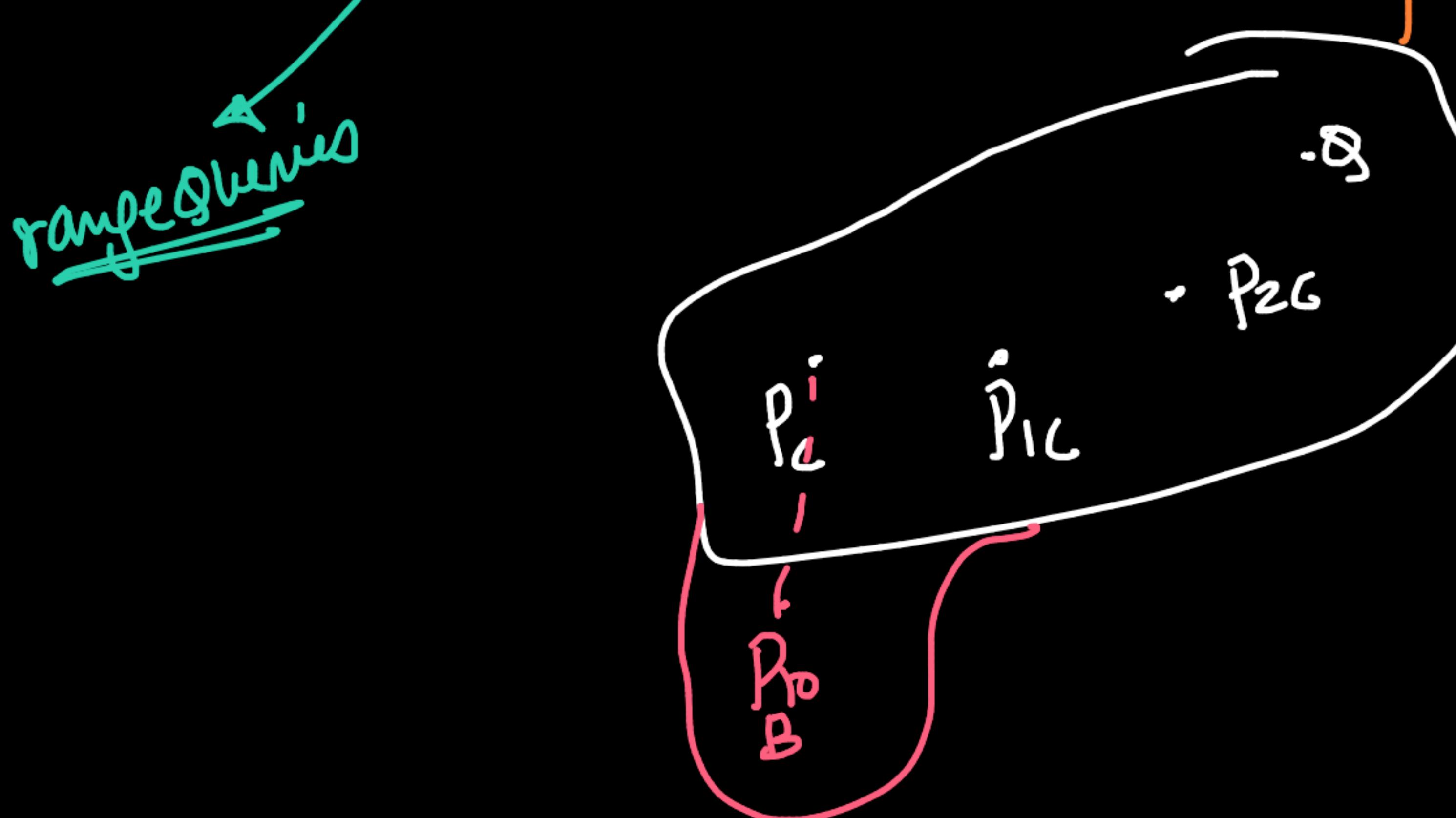
not assigned to any cluster

- create a new cluster with P
- add all points that are density connected to P in pts cluster

core-pt  
density connected  
range queries



plan  
④ + border pt ; assign it to the nearest  
co-e-pt's Cluster

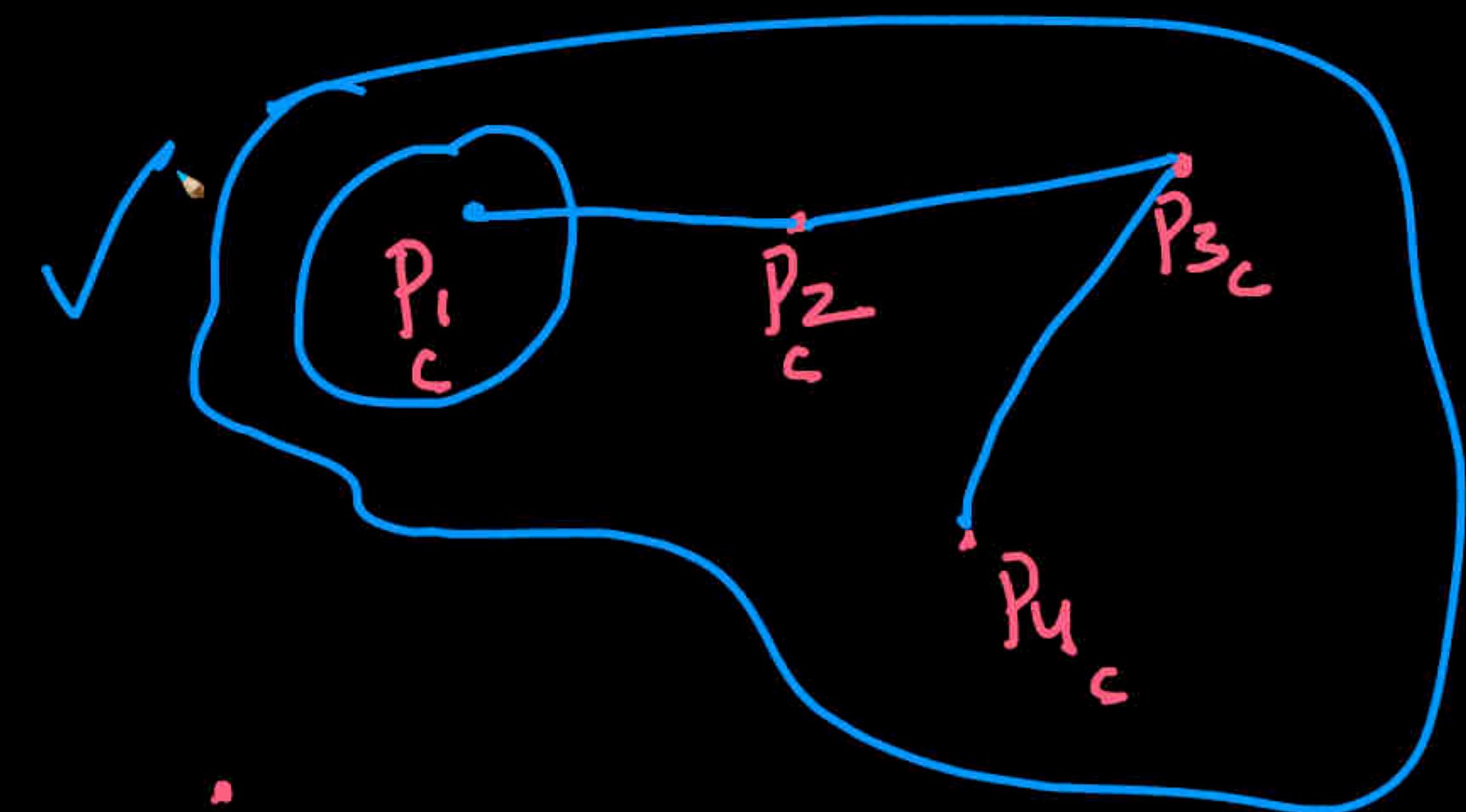


Time-complex:

$$\underline{\underline{O(n \lg n)}}$$

worst-case

- not noise sensitive
- computationally cheaper than Agg-clustering



P<sub>5B</sub>



~~Minpts~~

①

rules of thumb:



$$\text{Minpts} > d+1$$

$$\text{Typically: } - \text{Minpts} \approx 2 \times d$$

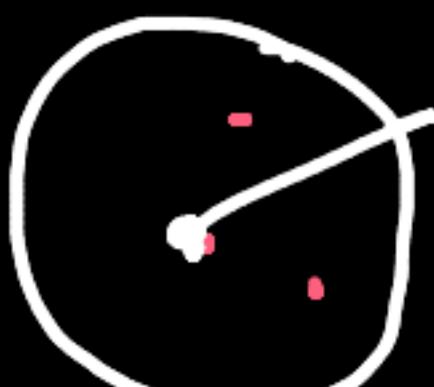
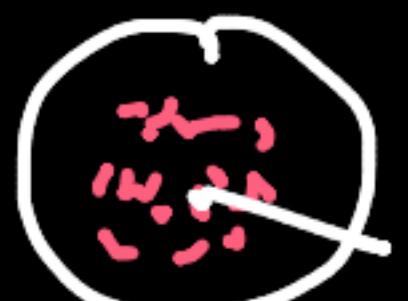
$$\frac{d}{2}$$

②

dataset

is noisy; fixed  $\epsilon_{\text{PS}}$

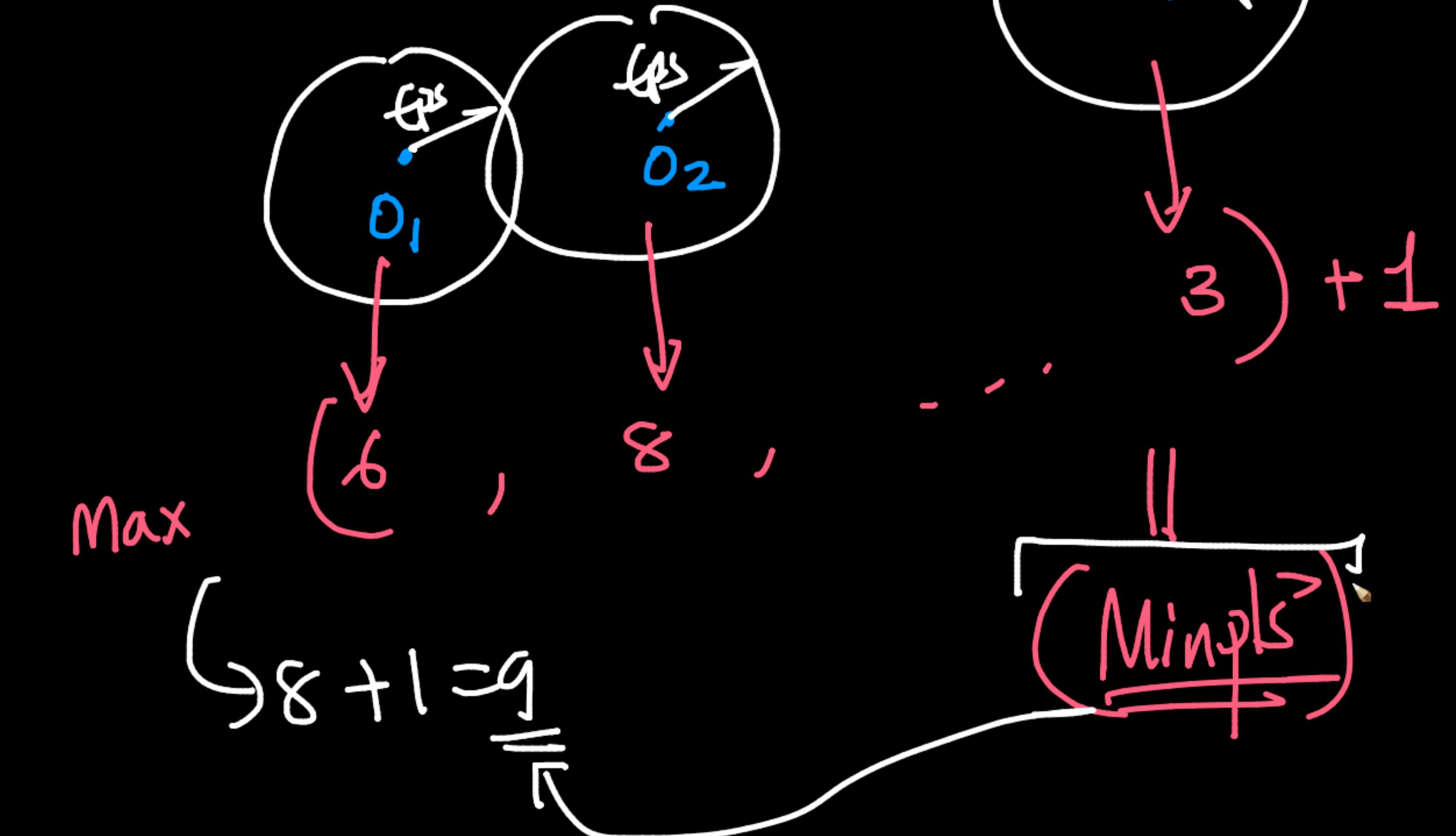
larger ~~Minpts~~



C

(hack)

Outliers using other methods: KNN



④

Various values of  
Minpts



different clusters



domain expert



Validation

~~Eps:~~~~step 1:~~

Minpts = 4  
(let)

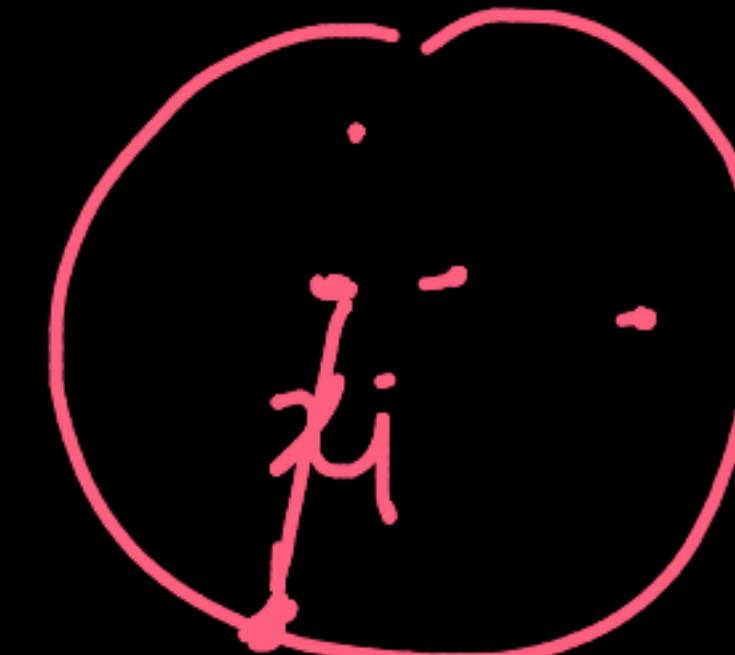
$$x_1 \rightarrow d_1$$

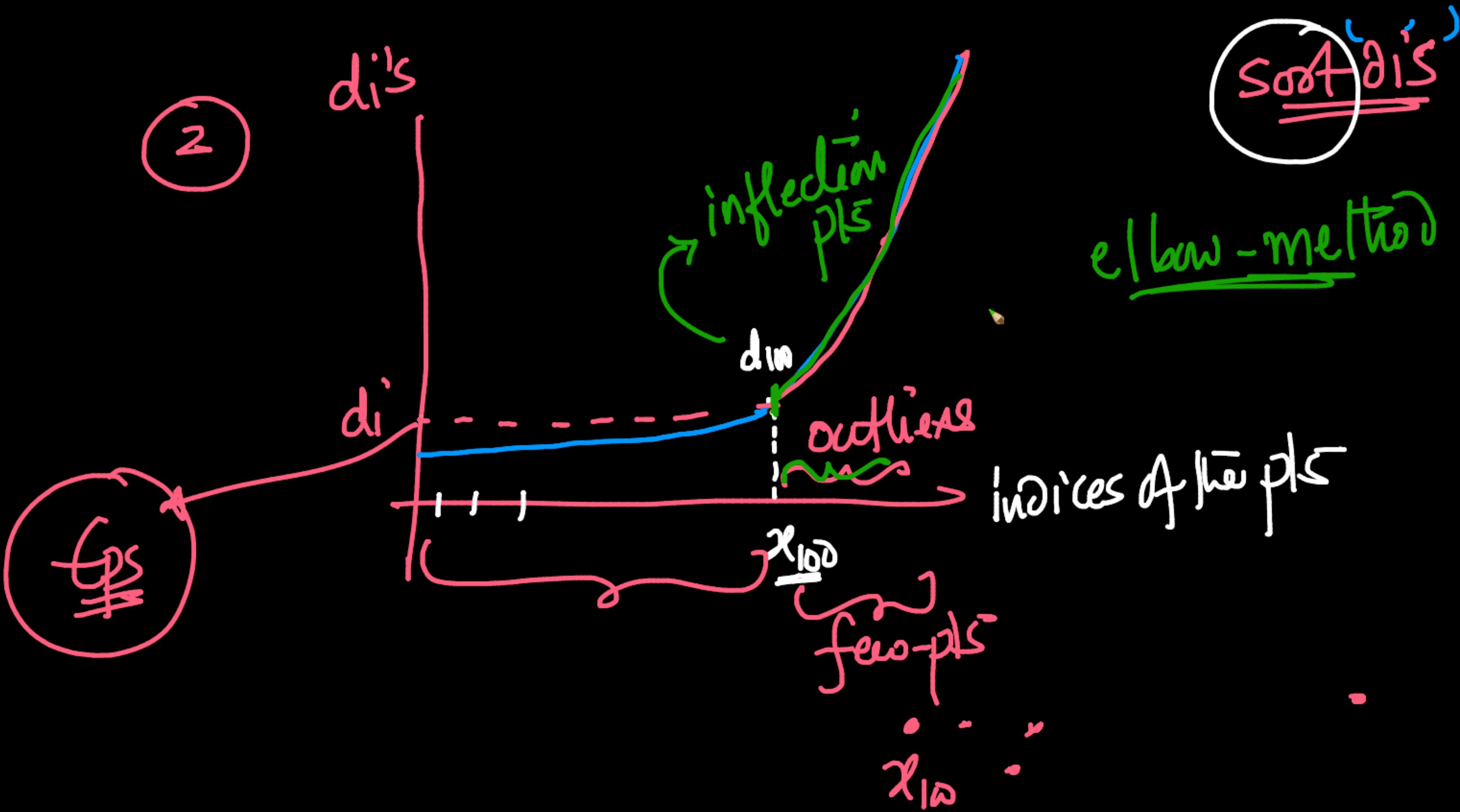
$$x_2 \rightarrow d_2$$

:

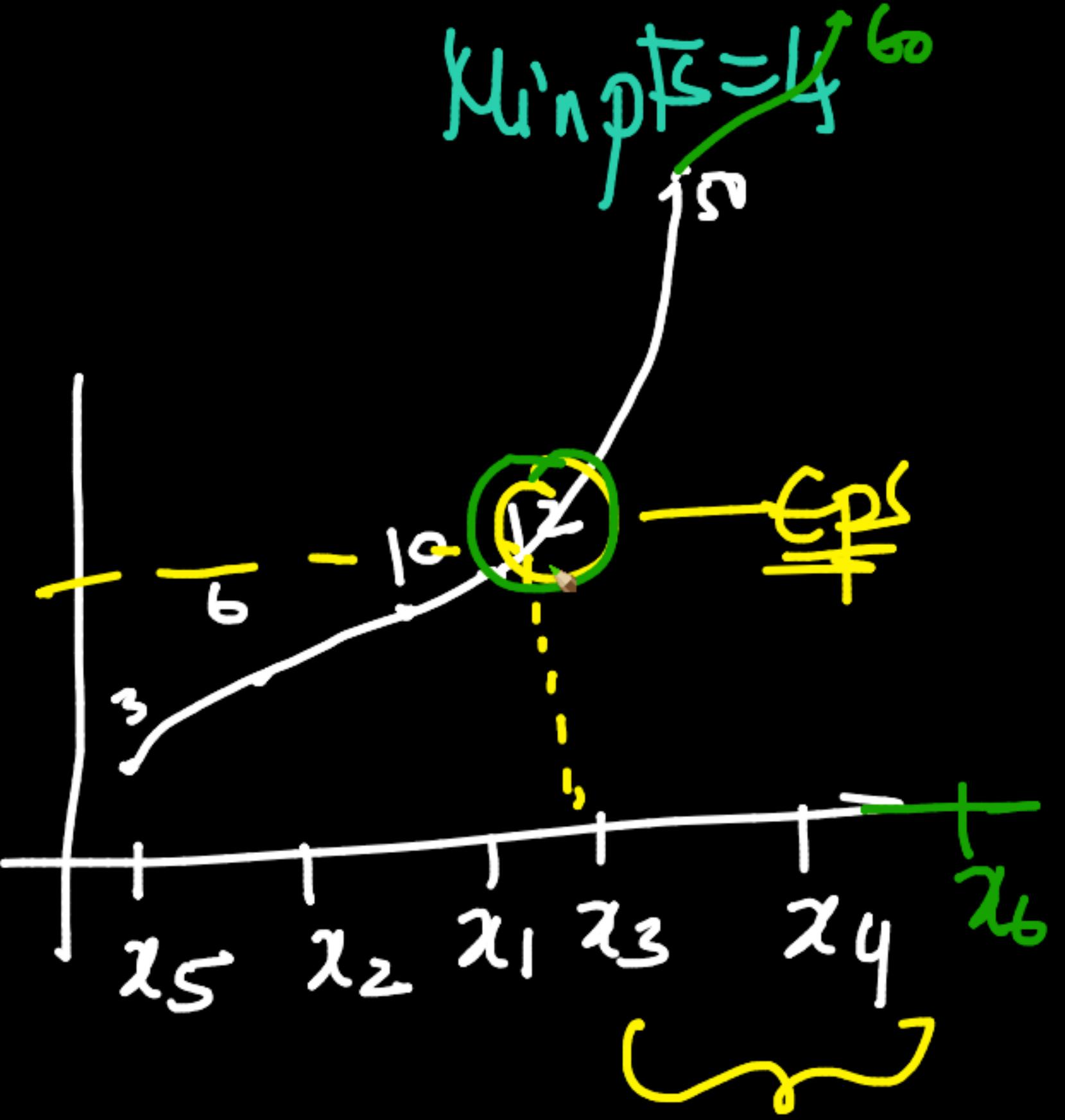
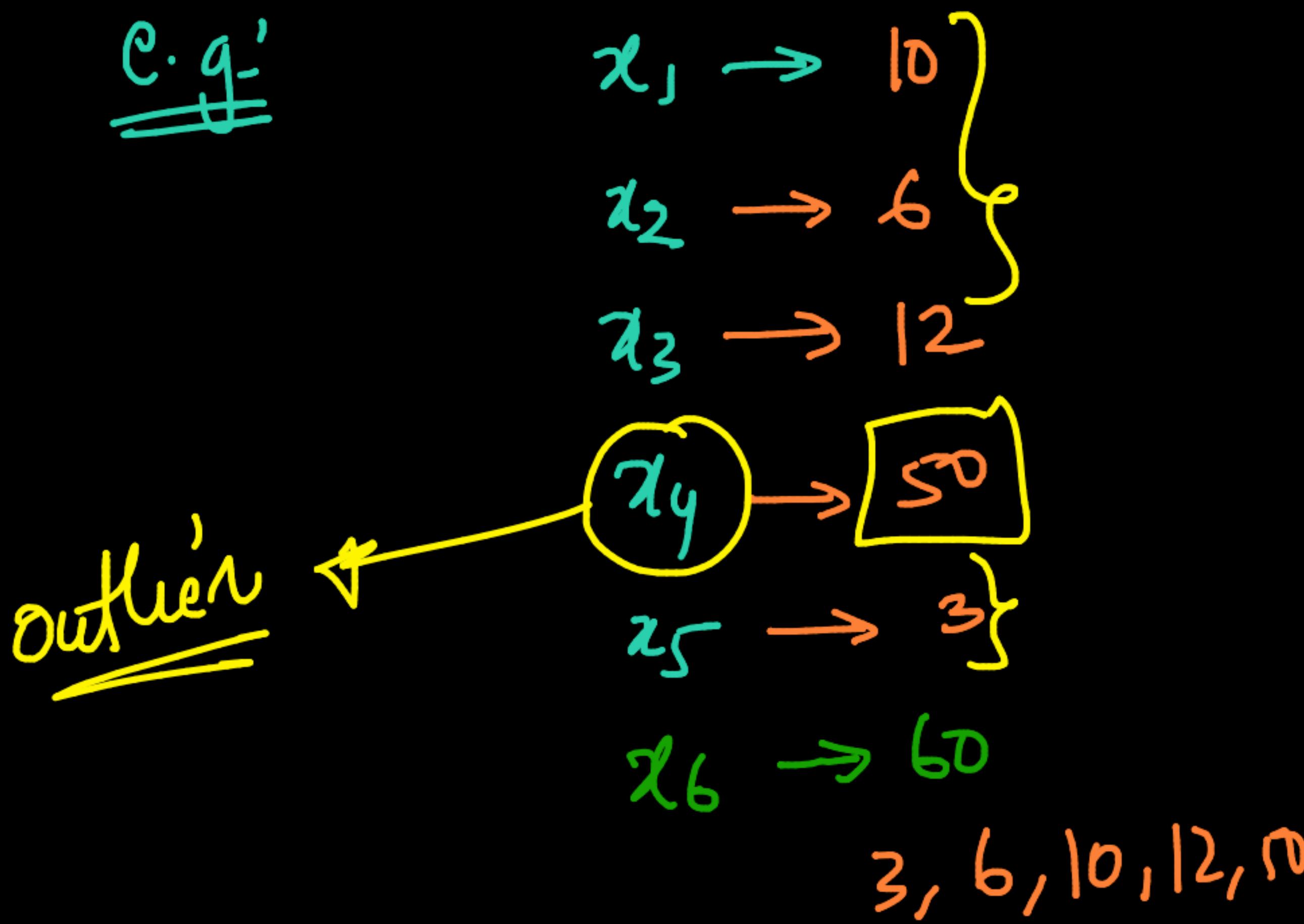
$\checkmark \{ x_i \rightarrow \underline{d_i} \rightarrow \text{distance from } x_i \text{ to } x_i \text{'s 4th NN}$

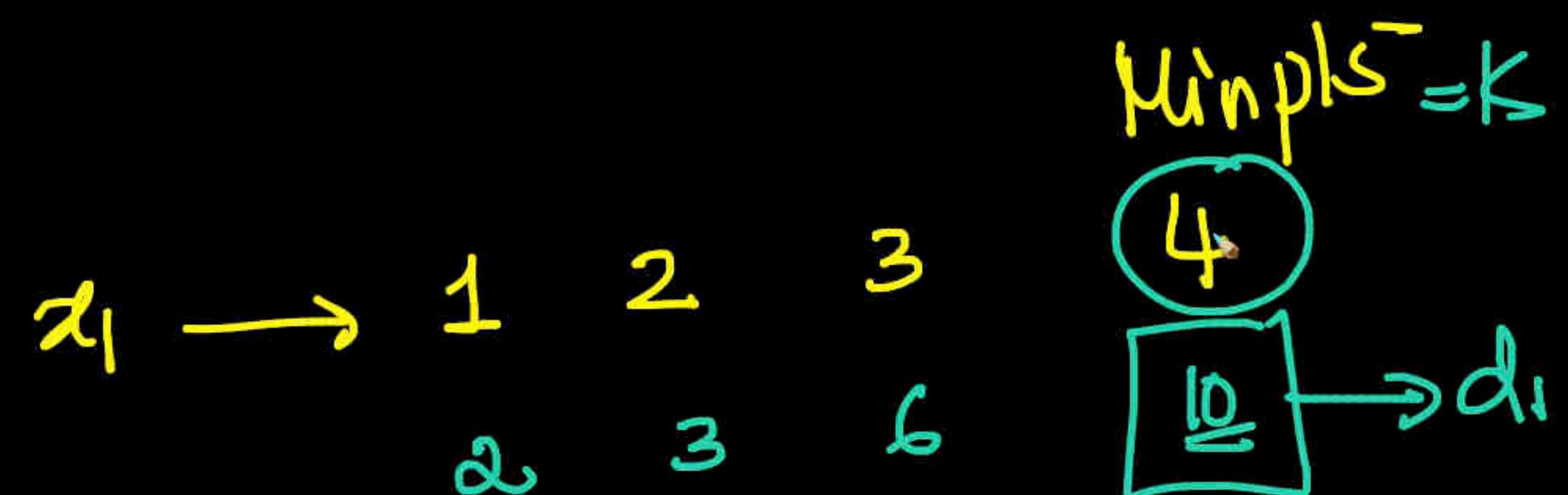
$$x_n \rightarrow d_n$$





C. g'

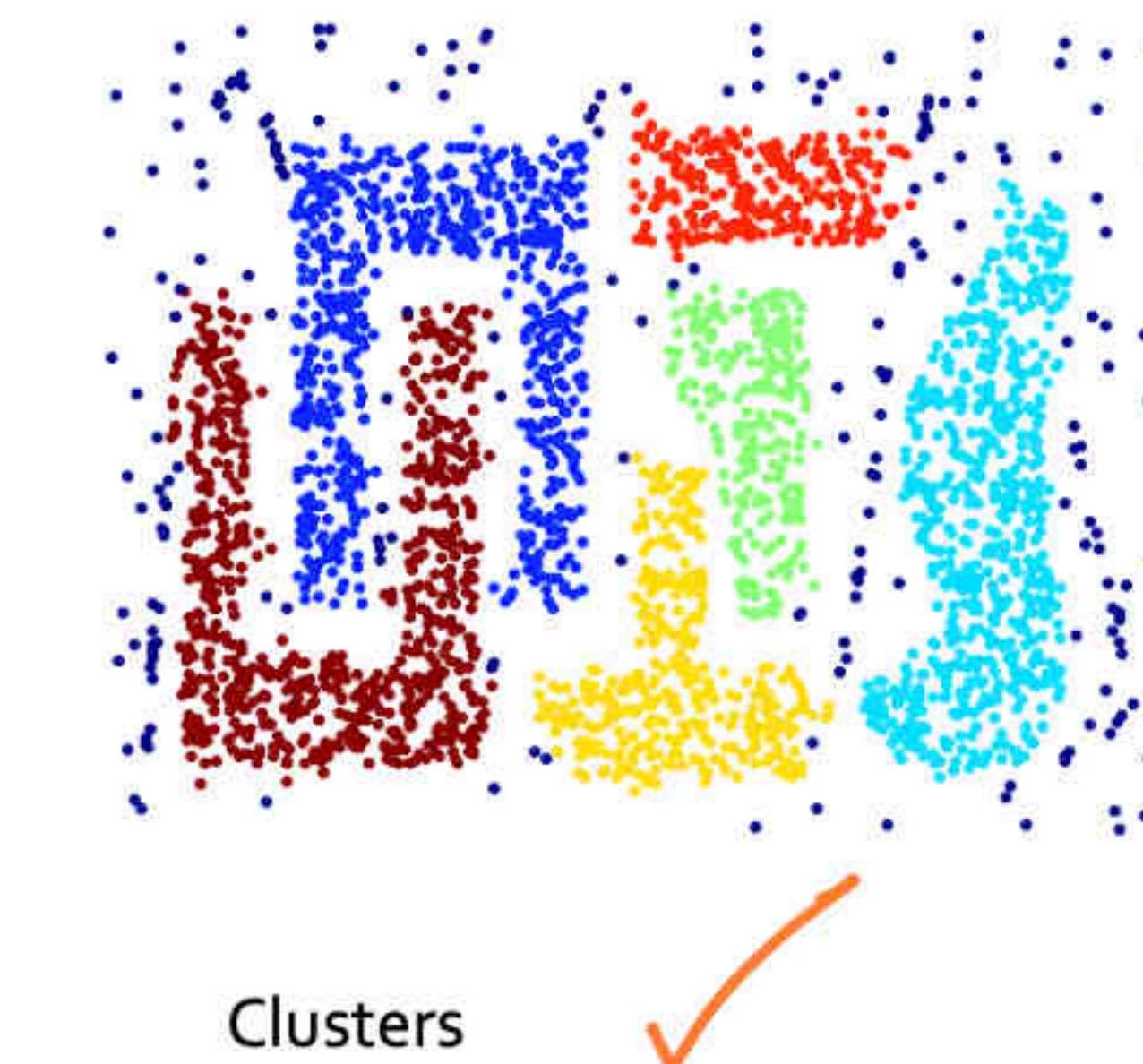
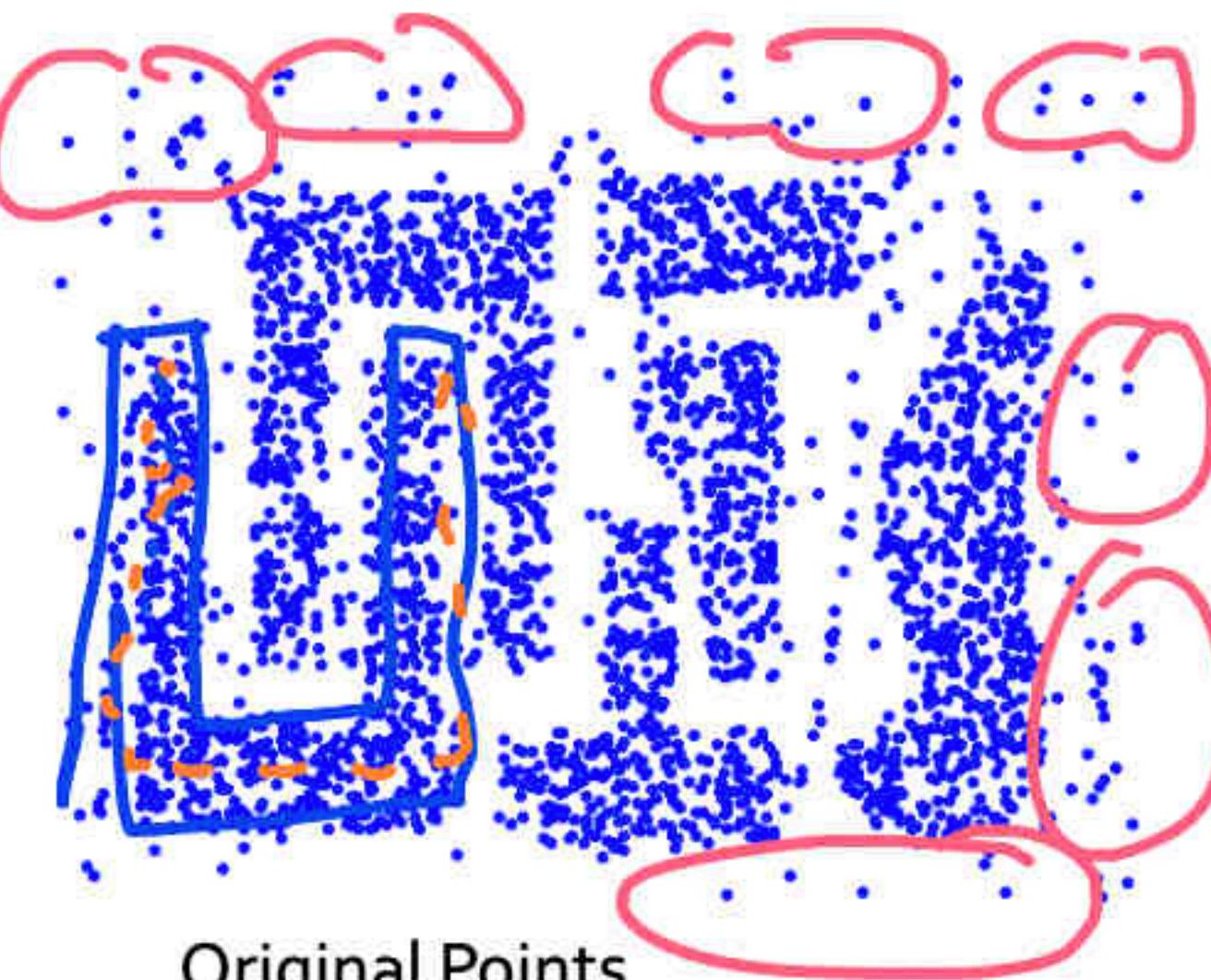




Mlnpts & Eps: ~~hyper-param~~ clusters → domain expert

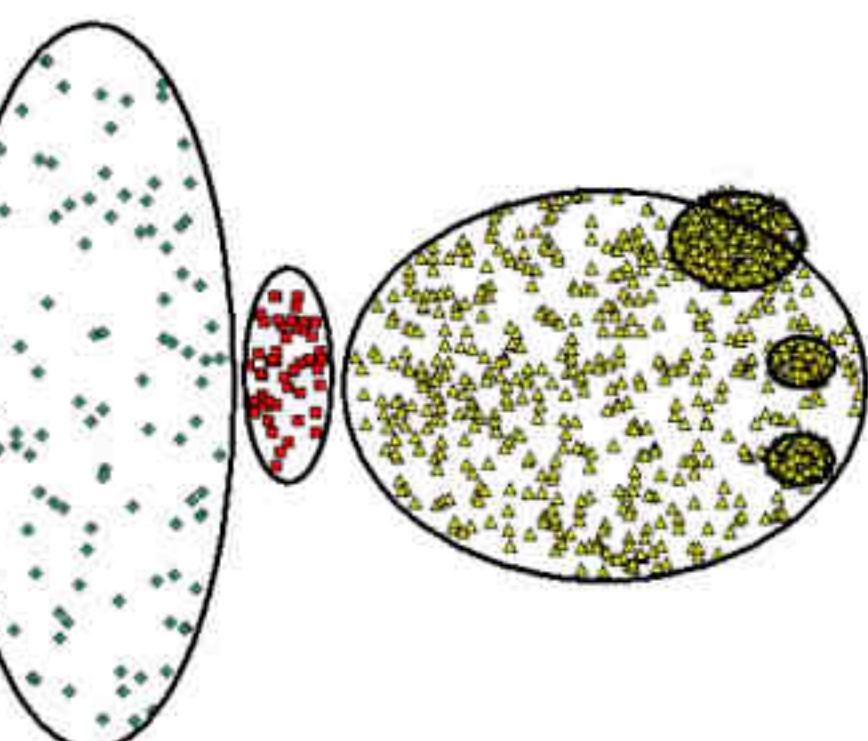
↓  
logical hacks to arrive @ reasonable values  
(heuristics)

# When DBSCAN Works Well



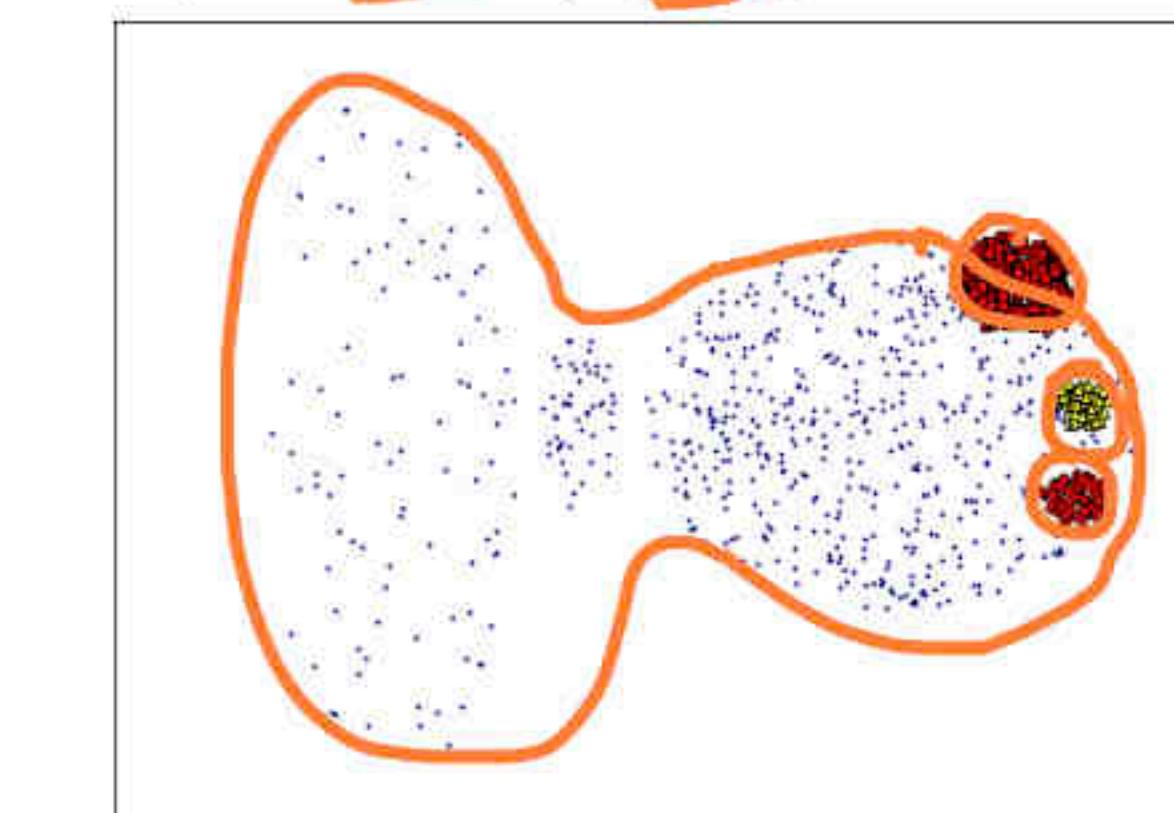
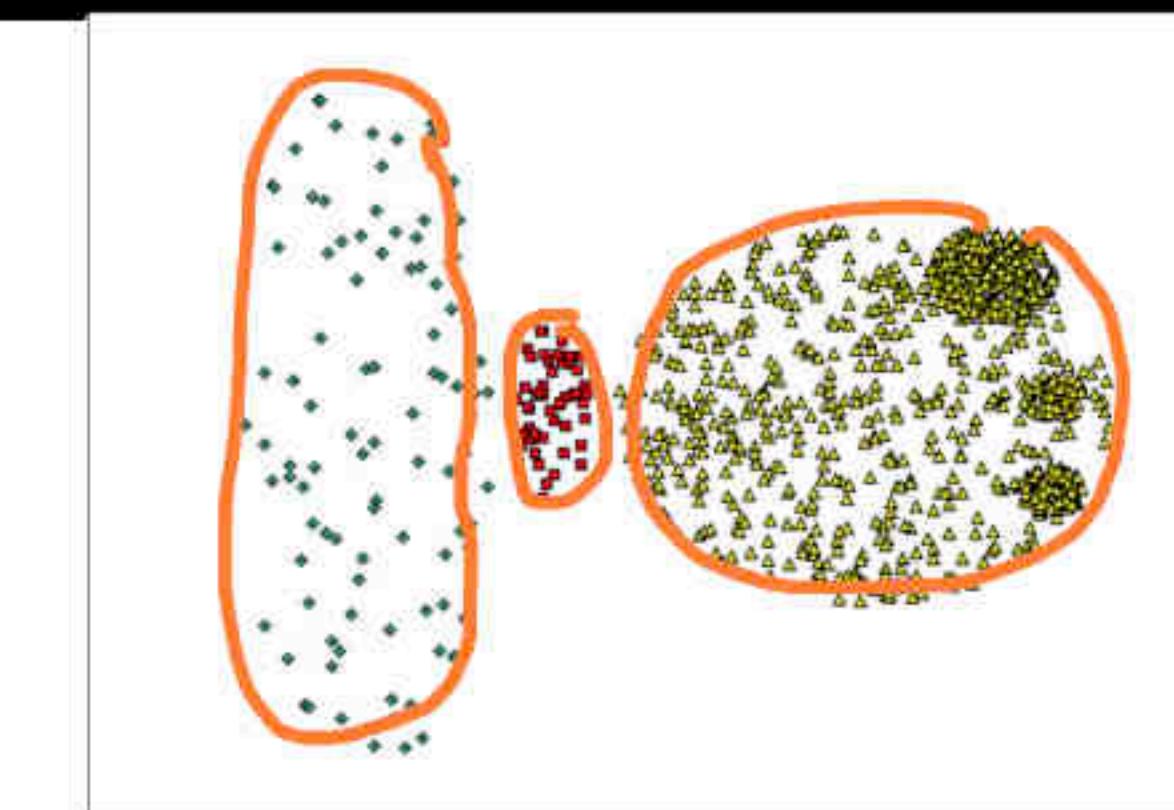
- Resistant to Noise
- Can handle clusters of different shapes and sizes

# Well



Original Points

- Varying densities
- High-dimensional data



sensitive to  
parameters:  
~~MinPts & Eps~~

# DBSCAN: Sensitive to Parameters $\leftarrow \text{0.5} \rightarrow \epsilon \rightarrow \leftarrow \text{0.4} \rightarrow$

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

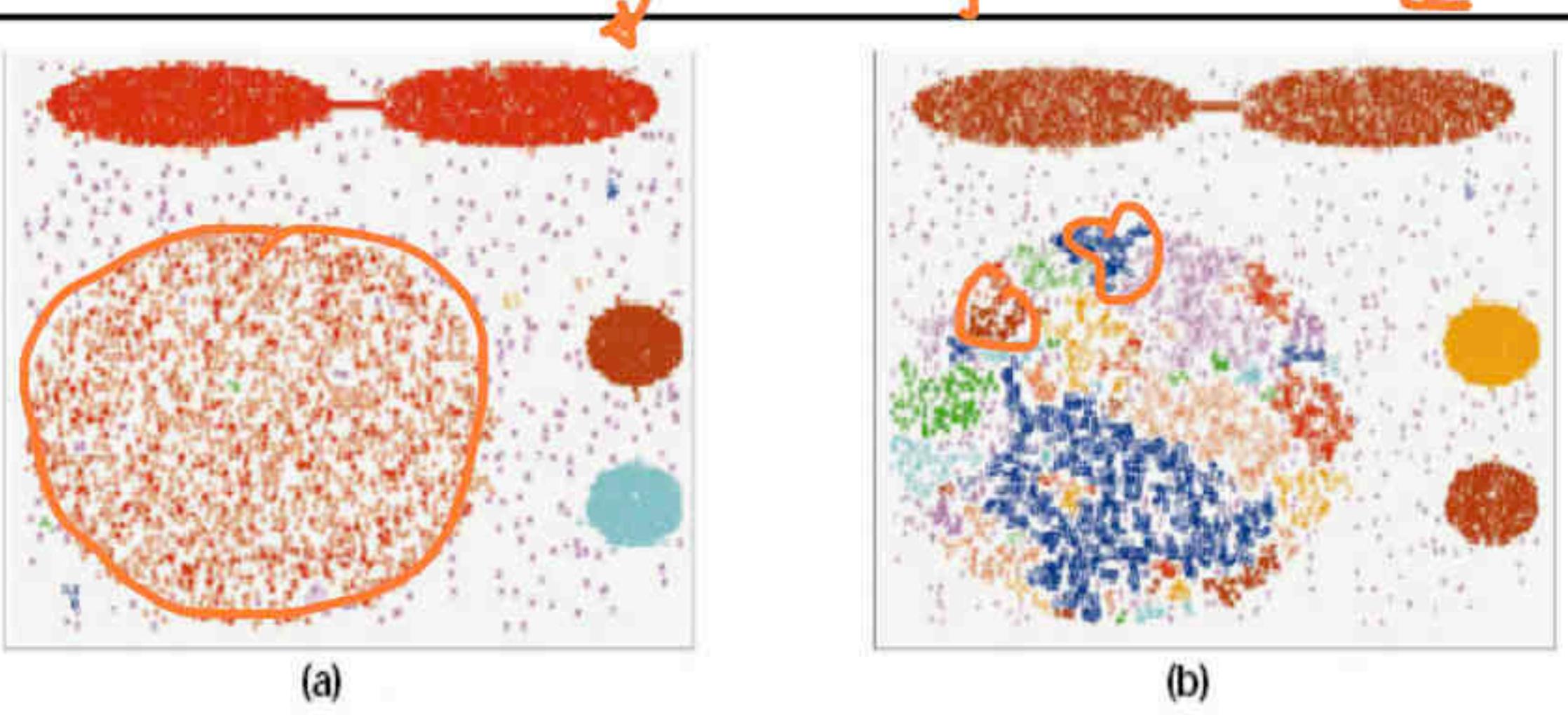
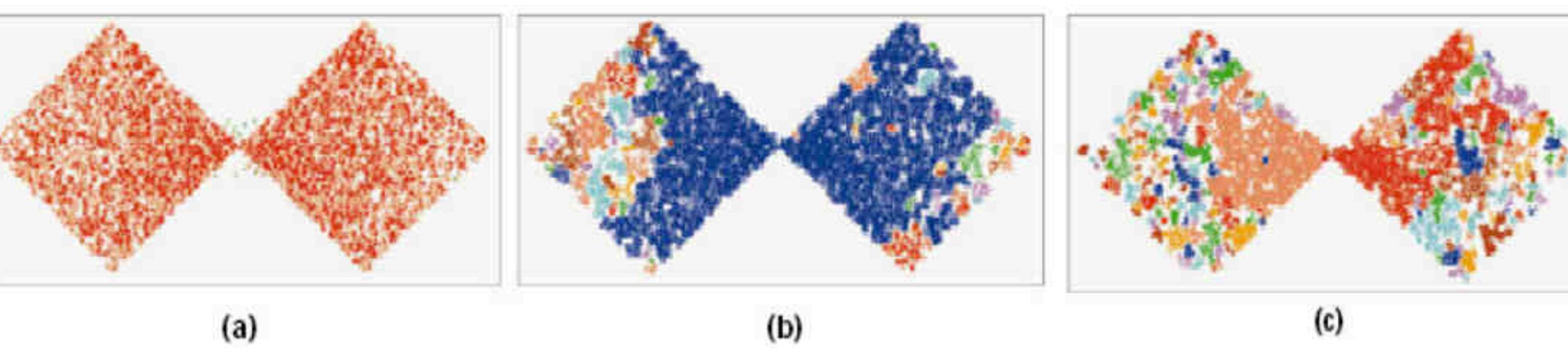


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



# DBSCAN: Sensitive to Parameters

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

2\* $\delta$

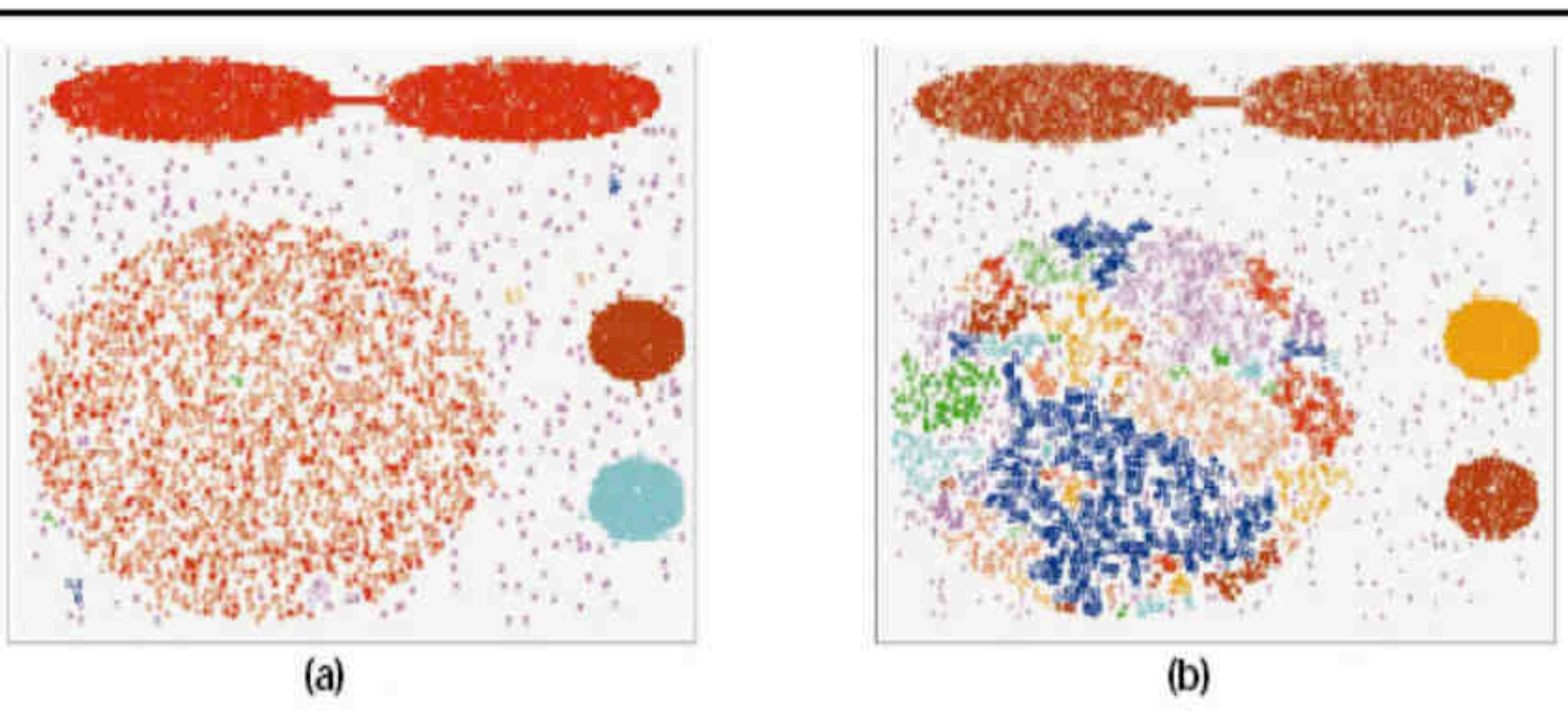
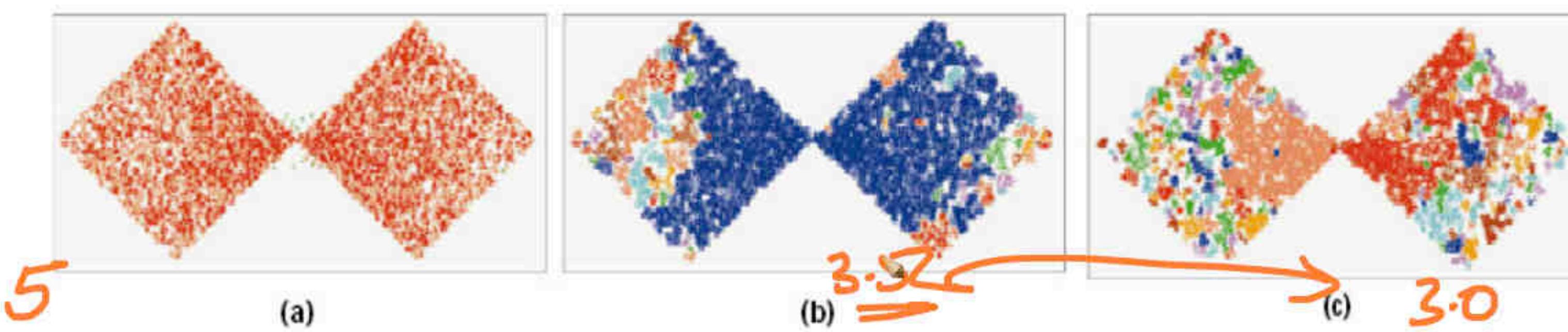


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



## DBScan

→ # clusters not needed like K-means

### Adv

→ non-globular

→ robust to outliers

→ databases because of efficient range queries  
(large data)

→ very sensitive to parameter: Minpts & eps  
disadv!

Microsoft PowerPoint - Cluster ... | DBSCAN - Wikipedia | sklearn.cluster.AgglomerativeClustering | sklearn.cluster.DBSCAN - scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html

scikit learn

Install User Guide API Examples Community More

Go

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.cluster.DBSCAN](#)

Examples using

[sklearn.cluster.DBSCAN](#)

## sklearn.cluster.DBSCAN

```
class sklearn.cluster.DBSCAN(eps=0.5, *, min_samples=5, metric='euclidean', metric_params=None,  
algorithm='auto', leaf_size=30, p=None, n_jobs=None)
```

[source]

Perform DBSCAN clustering from vector array or distance matrix.

DBSCAN - Density-Based Spatial Clustering of Applications with Noise. Finds core samples of high density and expands clusters from them. Good for data which contains clusters of similar density.

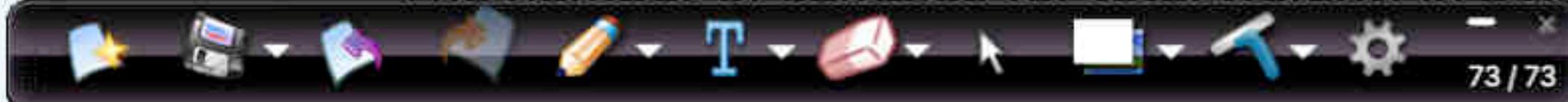
Read more in the [User Guide](#).

**Parameters:** *eps : float, default=0.5*

The maximum distance between two samples for one to be considered as in the neighborhood of the other. This is not a maximum bound on the distances of points within a cluster. This is the most important DBSCAN parameter to choose appropriately for your data set and distance function.

*min\_samples : int, default=5*

The number of samples (or total weight) in a neighborhood for a point to be point itself.



Microsoft PowerPoint - Cluster ... | DBSCAN - Wikipedia | sklearn.cluster.AgglomerativeClustering | sklearn.cluster.DBSCAN - scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html

scikit learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.cluster.DBSCAN

Examples using sklearn.cluster.DBSCAN

metric : string or callable, default='euclidean'  
metric is a string or callable, it must be one of the options allowed by `sklearn.metrics.pairwise_distances` for its metric parameter. If metric is "precomputed", X is assumed to be a distance matrix and must be square. X may be a sparse graph, in which case only "nonzero" elements may be considered neighbors for DBSCAN.

New in version 0.17: metric precomputed to accept precomputed sparse matrix.

metric\_params : dict, default=None  
Additional keyword arguments for the metric function.

New in version 0.19.

algorithm : {'auto', 'ball\_tree', 'kd\_tree', 'brute'}, default='auto'  
*R' lies* The algorithm to be used by the NearestNeighbors module to compute pointwise distances and find nearest neighbors. See NearestNeighbors module documentation for details.

leaf\_size : int, default=30  
Leaf size passed to BallTree or cKDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.

p : float, default=None  
None, then p=2 (equivalent to the Euclidean distance).

Microsoft PowerPoint - Cluster ... | DBSCAN - Wikipedia | sklearn.cluster.AgglomerativeClustering | sklearn.cluster.DBSCAN - scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html

scikit learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.cluster.DBSCAN

Examples using sklearn.cluster.DBSCAN

algorithm : {'auto', 'ball\_tree', 'kd\_tree', 'brute'}, default='auto'

The algorithm to be used by the NearestNeighbors module to compute pointwise distances and find nearest neighbors. See NearestNeighbors module documentation for details.

leaf\_size : int, default=30

Leaf size passed to BallTree or cKDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.

p : float, default=None

The power of the Minkowski metric to be used to calculate distance between points. If None, then p=2 (equivalent to the Euclidean distance).

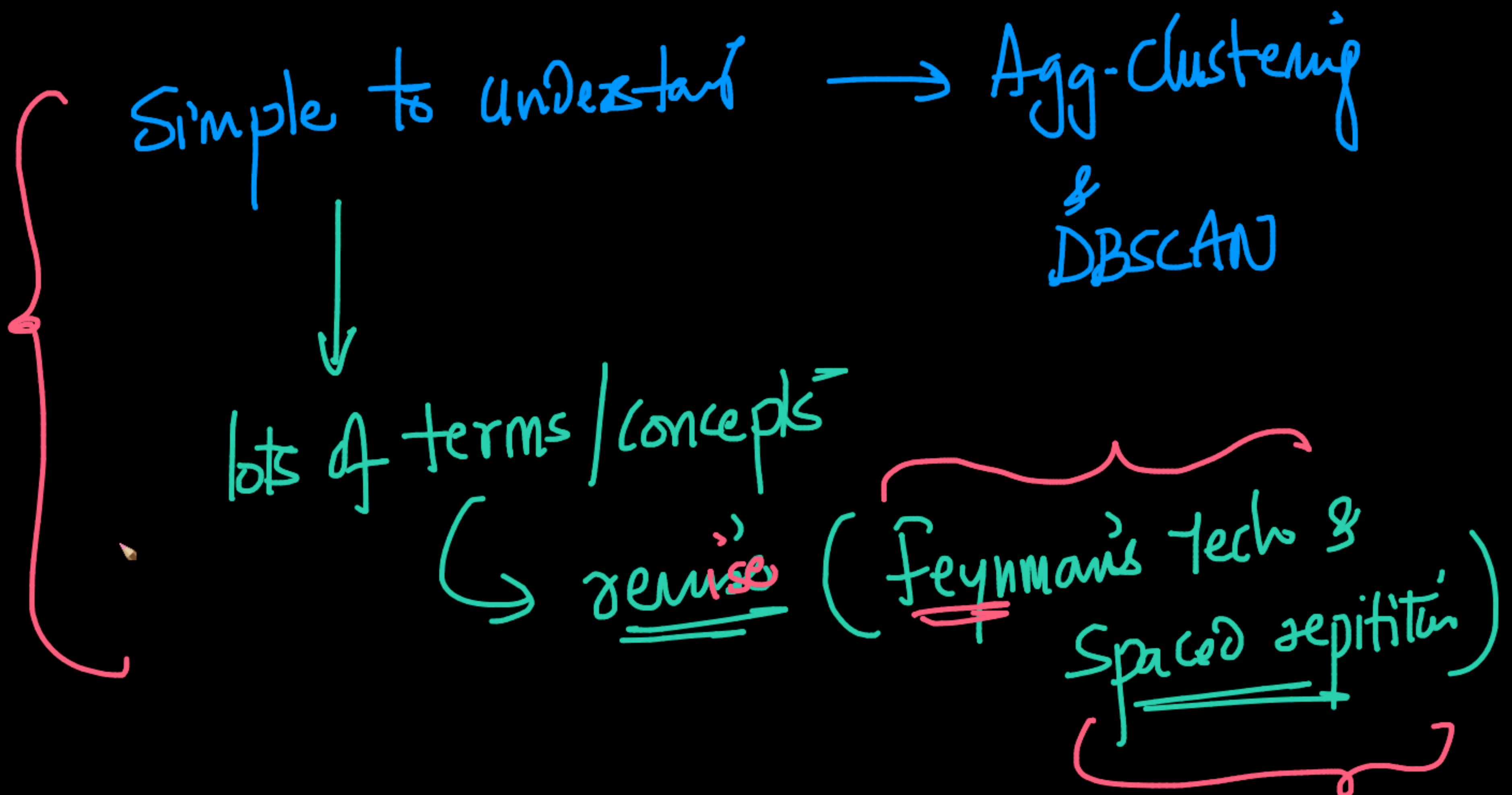
n\_jobs : int, default=None

The number of parallel jobs to run. None means 1 unless in a joblib.parallel\_backend context. -1 means using all processors. See Glossary for more details.

Attributes:

- core\_sample\_indices\_ : ndarray of shape (n\_core\_samples,) Indices of core samples.
- components\_ : ndarray of shape (n\_core\_samples, n\_features)

Toggle Menu





Prev Up Next

scikit-learn 1.1.1

[Other versions](#)Please [cite us](#) if you use the software.[sklearn.cluster.DBSCAN](#)

Examples using

[sklearn.cluster.DBSCAN](#)

# sklearn.cluster.DBSCAN

```
class sklearn.cluster.DBSCAN(eps=0.5, *, min_samples=5, metric='euclidean', metric_params=None,  
algorithm='auto', leaf_size=30, p=None, n_jobs=None)
```

[\[source\]](#)

Perform DBSCAN clustering from vector array or distance matrix.

DBSCAN - Density-Based Spatial Clustering of Applications with Noise. Finds core samples of high density and expands clusters from them. Good for data which contains clusters of similar density.

Read more in the [User Guide](#).

**Parameters:** **eps : float, default=0.5**

The maximum distance between two samples for one to be considered as in the neighborhood of the other. This is not a maximum bound on the distances of points within a cluster. This is the most important DBSCAN parameter to choose appropriately for your data set and distance function.

**min\_samples : int, default=5**



neighborhood for a point to be

considered as a core point. This includes the point itself.