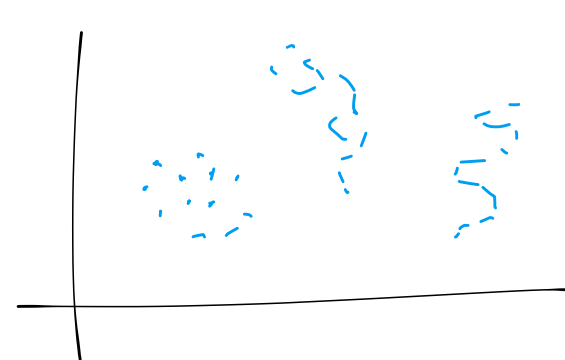
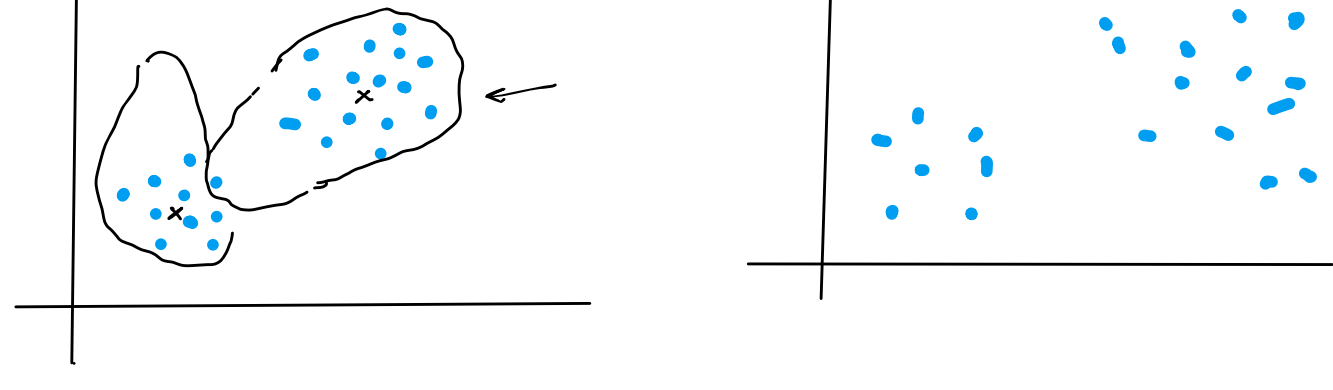


## Clustering

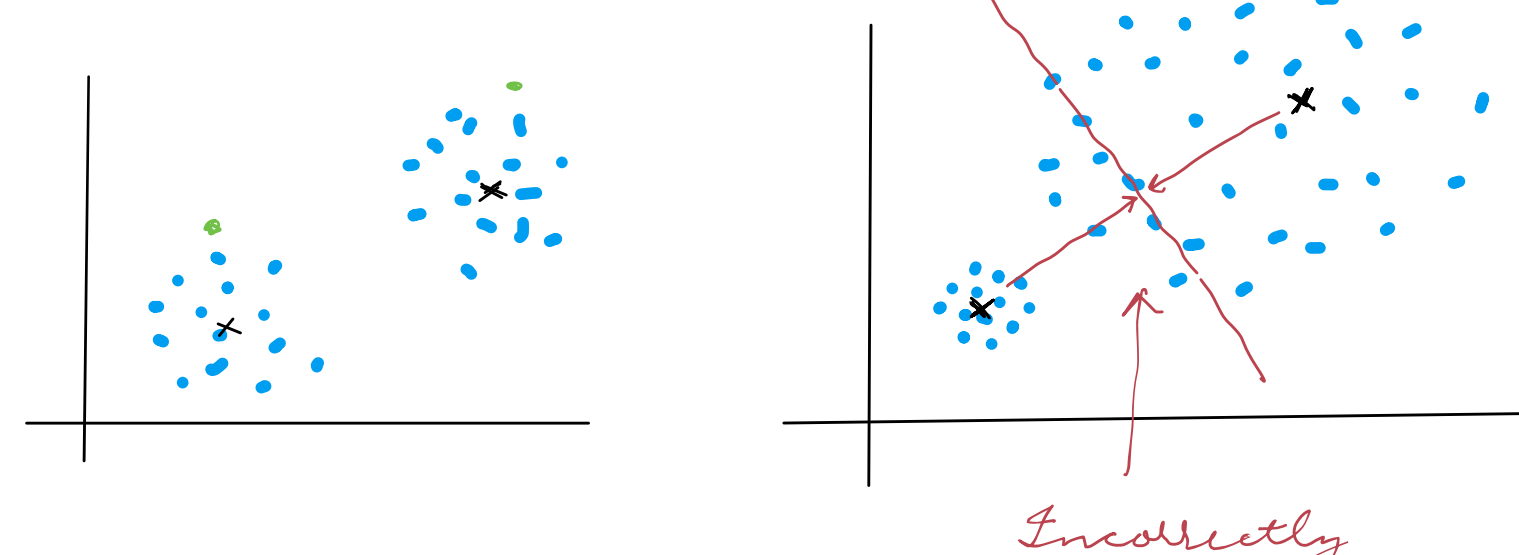


No labels  
 Supervised  $\rightarrow X, Y$   
 $\hookrightarrow X \rightarrow Y$

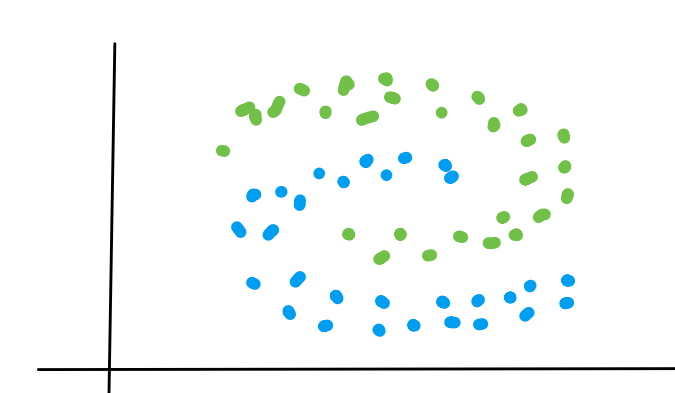
## K-Means clustering



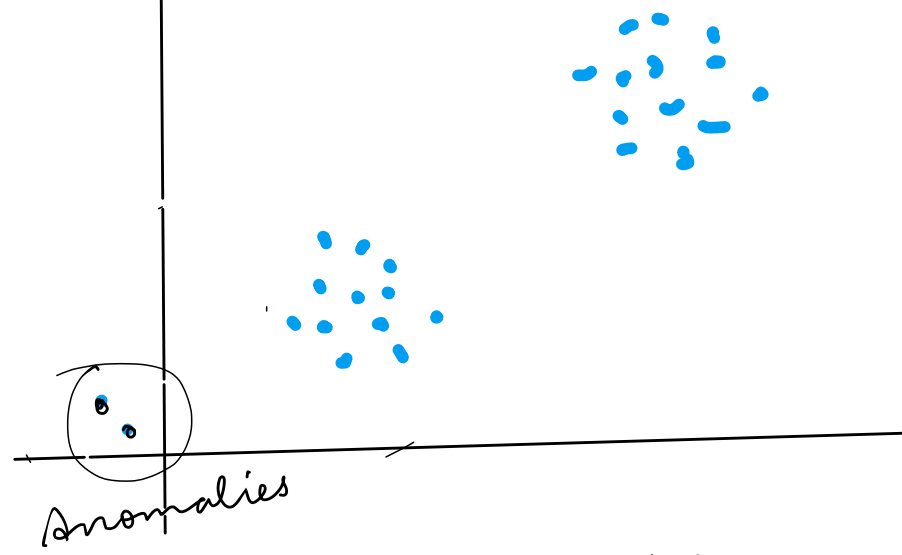
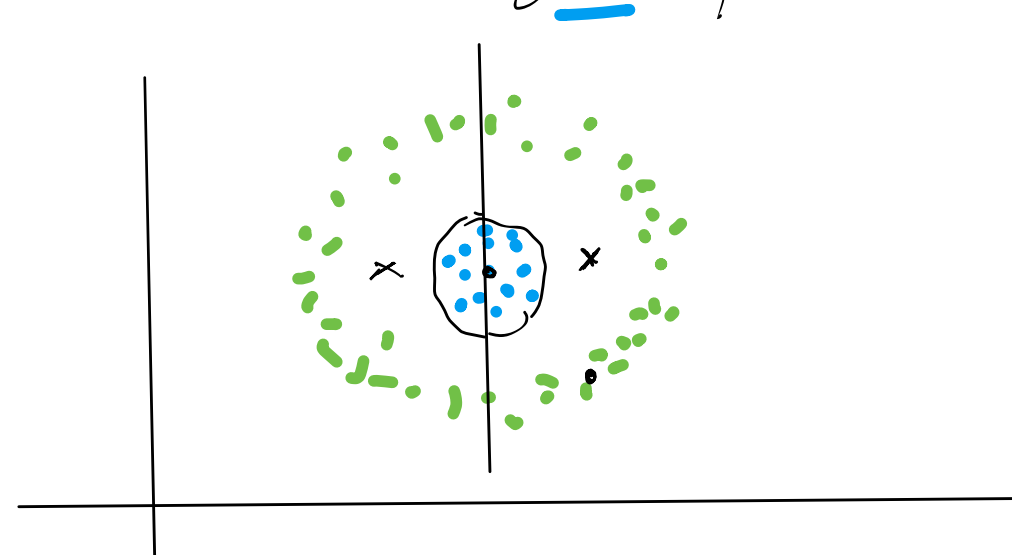
1. K clusters
2. Centroids  $\leftarrow$  Initialize  
 $\hookrightarrow$  Take the average of each dimension for all the points in the cluster.
4. Recompute Centroids  
 $\hookrightarrow$  Mean
3. Compute the distance of each point to each centroid. Assign it to the closest
5. Repeat till centroids stop updating or a preset number of iterations have passed.



*Incorrectly clustered*  
 • Cannot handle varying densities



- Creates spherical clusters
- Creates equal partitions

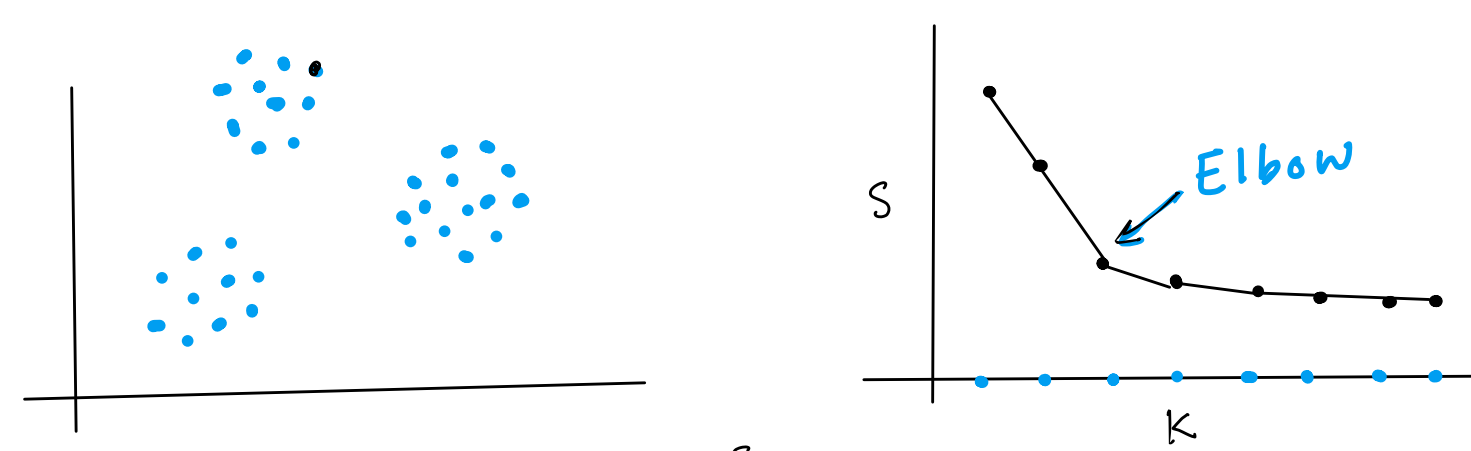


Non-core  
 $\hookrightarrow$  Non-DR

How do you select K?

$\hookrightarrow$  Elbow method

KMeans } Cluster vs Anomaly  
 HC }  
 DBSCAN  $\leftarrow$  Anomaly



K=1  $\rightarrow$  KMeans  $\rightarrow$  Score  
 K=2  $\rightarrow$  KMeans  $\rightarrow$  Score  
 K=3  $\rightarrow$   $\rightarrow$   $\rightarrow$   
 4  
 5  
 $\vdots$   
 10

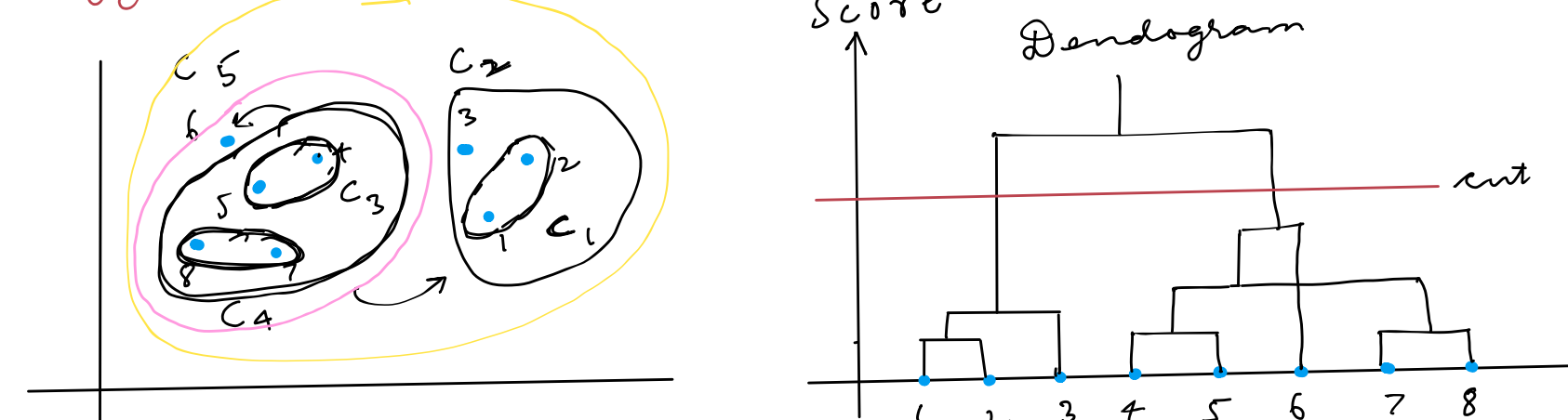
Silhouette score

a  $\rightarrow$  Average intra cluster dist  
 b  $\rightarrow$  Distance to closest point from another cluster

$$\frac{b-a}{\max(a,b)} \rightarrow 1$$

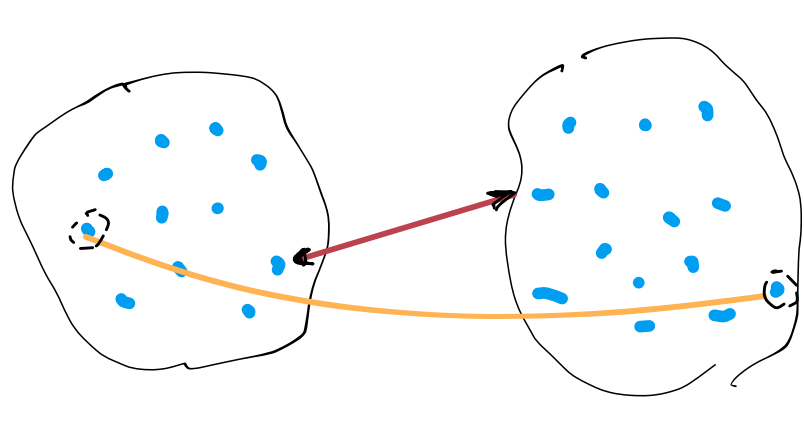
Avg. intra clust.  $\leftarrow$   
 Avg. inter clust.  $\rightarrow$

## Agglomerative Hierarchical clustering



1. Initialize each point as a separate cluster.
2. Repeat till 1 cluster remains  
 $\hookrightarrow$  Merge closest clusters

$\downarrow$  need  
 Distance metrics for clusters



1. Single linkage  
 $\hookrightarrow$  Distance between the closest 2 points in the clusters.
2. Complete linkage  
 $\hookrightarrow$  Distance between the 2 farthest points.  
 $\hookrightarrow$  Intuition: Diameter of the cluster if the smaller 2 are merged.
3. Average linkage  
 $\hookrightarrow$  Average the distance between all pairs of points in the two clusters.

Limitations

- $\hookrightarrow$  Where to cut the dendrogram?
- $\hookrightarrow$  Which linkage to use?
- $\hookrightarrow$  Expensive

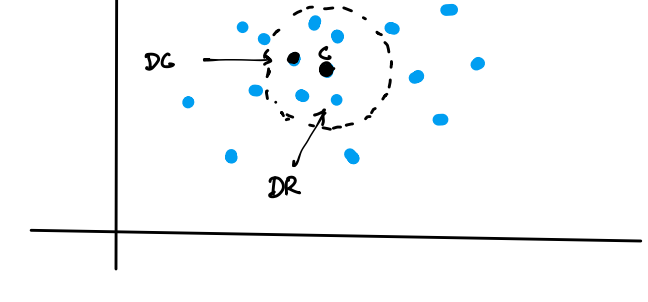
## DBSCAN

$\hookrightarrow$  You don't have to specify the #clusters

Terminology

$\epsilon$  : radius

minSep : #neighbours to label a point as core point.



If a point has more than minSep points in its  $\epsilon$ -neighborhood  
 $\hookrightarrow$  It is a core point.

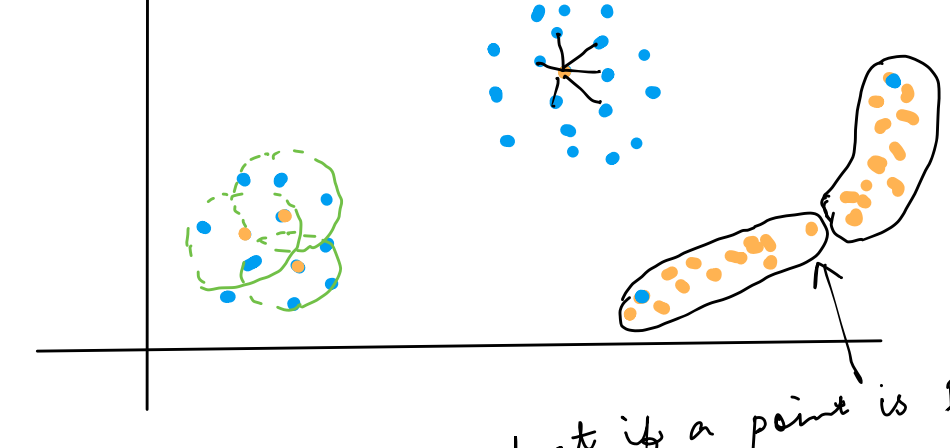
① Core points.

$\hookrightarrow$  Points inside the  $\epsilon$ -neighborhood of a core point are said to be

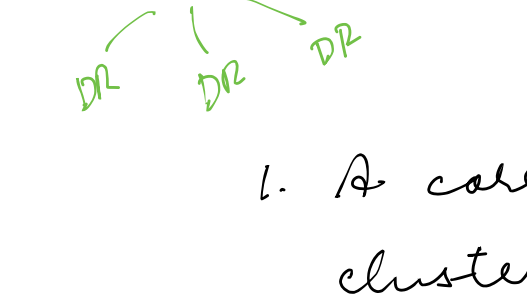
② "Density Reachable" (DR)

③ Density Connected (DC)

$\hookrightarrow$  If a "density reachable" point is itself core point, it is said to be "density connected."



What if a point is DR from 2 diff clusters?



1. A core point can start its own cluster.
2. DR points belong to their core point's clusters.
3. A core point expands its cluster through DC points.
4. Non-core, non DR points are anomalies.