# Lecture 22 - Apr 8

Density-based Clustering
DBSCAN
OPTICS

References

*Data Mining and Machine Learning*
Ch 15 - Density-based Clustering

Where K-Means and EM Fails...


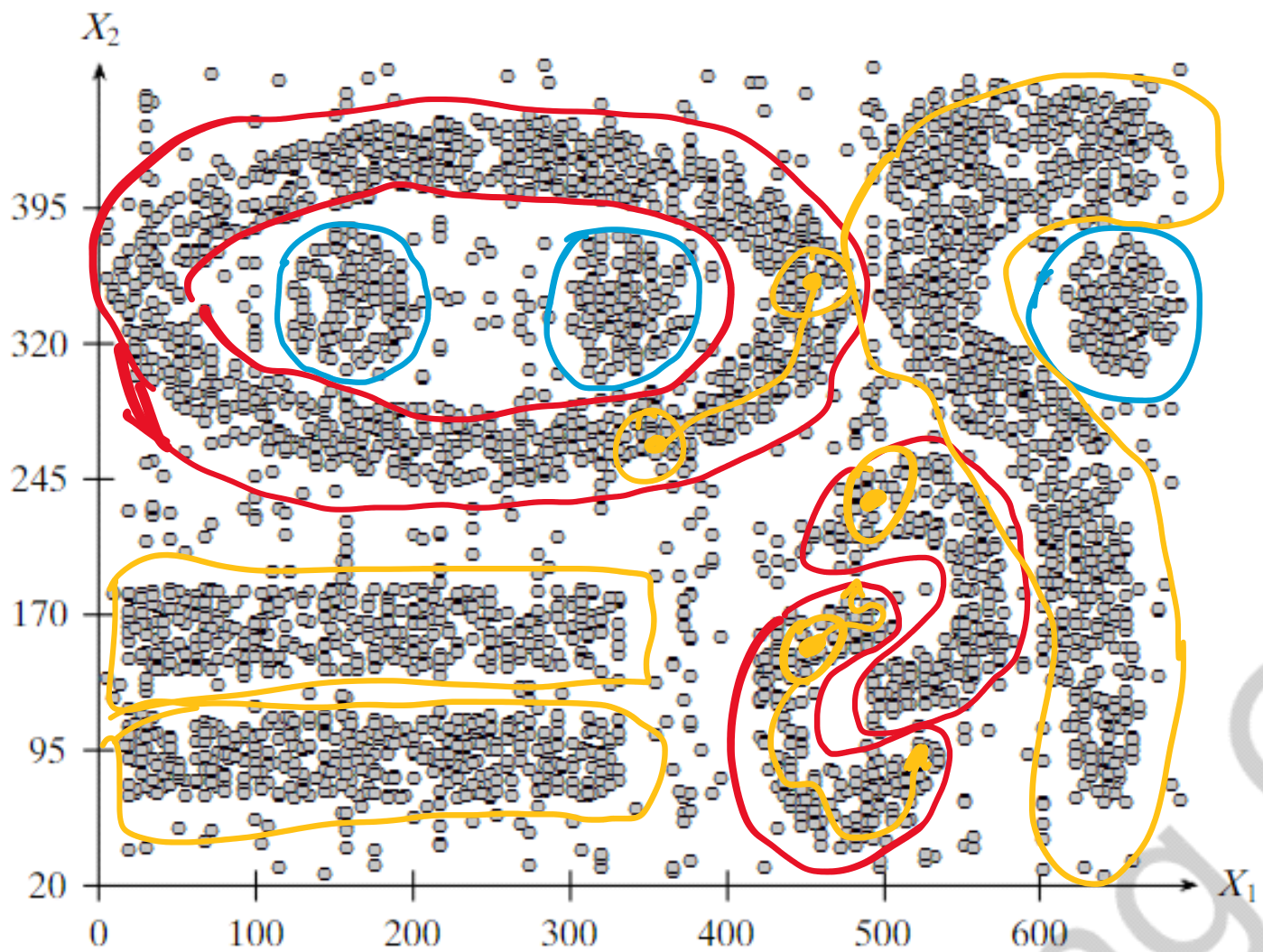
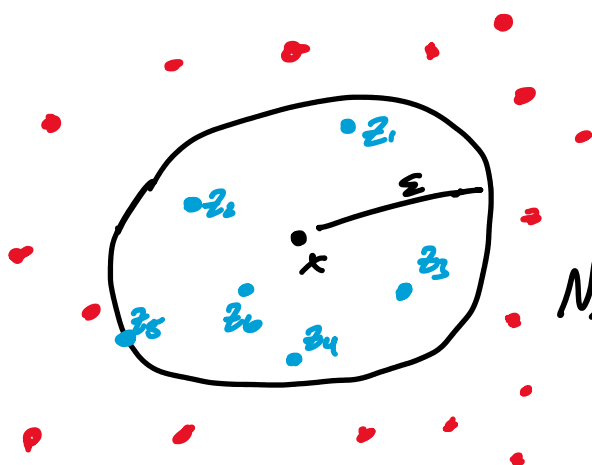**Figure 15.1.** Density-based dataset.

$X = \{x_1, \ldots, x_n\} = $ data points.

Density-based clustering uses local density of points to determine clusters rather than only distance between points.

We define a ball of radius $\varepsilon$ around a point $x \in \mathbb{R}^d$, called an $\varepsilon$-neighborhood of $x$ as

$$N_\varepsilon(x) = B_d(x, \varepsilon) = \left\{ z \in X \mid \underbrace{\|x - z\|_2}_{} \leq \varepsilon \right\}$$

<span style="color:blue">Euclidean distance, but others may be used</span>



$N_\varepsilon(x) = \{z_1, z_2, \ldots, z_6\}$

<span style="color:blue">user-defined local density threshold</span>

$x \in X$ is a <u>core point</u> if $|N_\varepsilon(x)| \geq m$

$x \in X$ is a <u>border point</u> if $|N_\varepsilon(x)| < m$ and it belongs to an $\varepsilon$-neighborhood of some core point
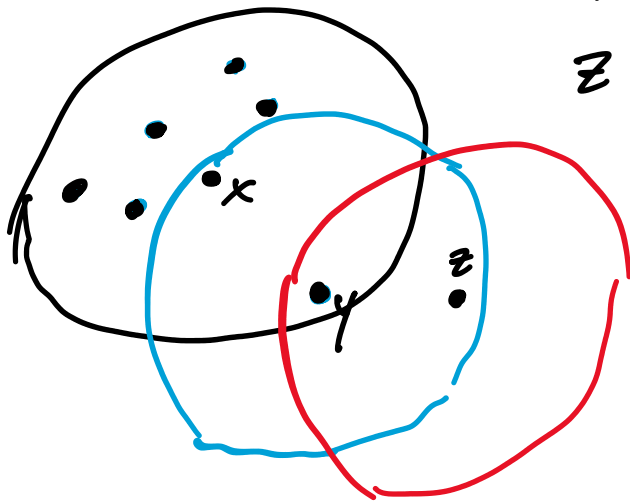
$\left[$ i.e. for some $z \in X$ with sufficiently dense neighborhood s.t. $x \in N_\varepsilon(z) \right]$

border point. it

If $x \in X$ is not a core point nor a border point, it is a noise point (or outlier)

Example: Let $m = 6 \Rightarrow$ $x$ is a core point

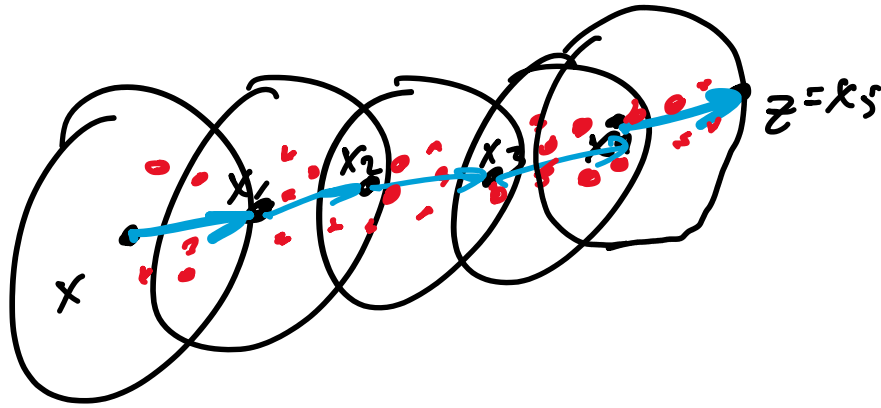$y$ is a border point

$z$ is a noise point



$x$ is directly density reachable from $y$ if $x \in N_\varepsilon(y)$ and $y$ is a core point

$x$ is density reachable from $z$ if there exists a chain of points $x_0, x_1, \ldots, x_\ell$ with $x = x_0$ and $z = x_\ell$ such that $x_i$ is directly density reachable from $x_{i-1}$, for all $i$.
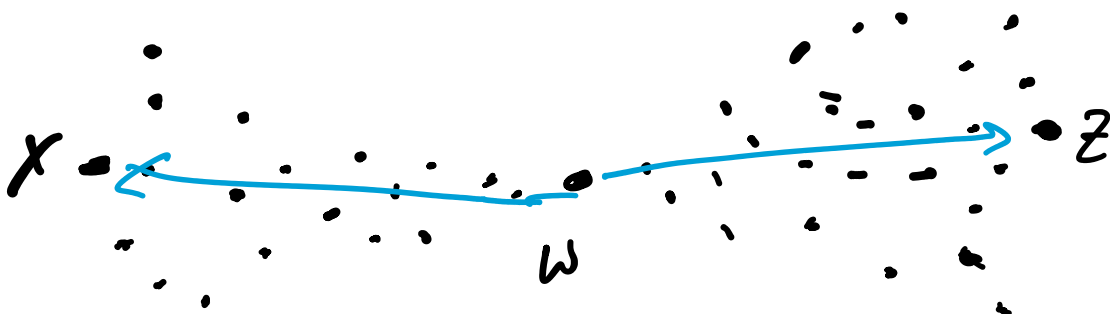
Such that $x_i$ is *collectivly density* ...

↳ i.e there is a set of core points leading from $z$ to $x$.



Density reachability is asymmetric, i.e. it is possible for $x$ to be D.R from $z$ but $z$ not be D.R. from $x$.

$x$ and $z$ are <u>density connected</u> if there exists a point $w$ s.t. $x$ and $z$ are both D.R from $w$.

A density-based cluster is a maximal set of
density connected points
$\hookrightarrow$ note it can be any shape, not just convex shapes

# DBSCAN Algorithm

inputs: $X, \varepsilon, m$

$\varepsilon$: radius of neighbourhood

$m$: local density threshold

$Core \leftarrow \emptyset$

for $x_i \in X$:

    Compute $N_\varepsilon(x_i)$

    $id(x_i) \leftarrow \emptyset$   ← set cluster id to $\emptyset$

    if $|N_\varepsilon(x_i)| \geq m$:   $Core \leftarrow Core \cup \{x_i\}$

> Compute $\varepsilon$-neighbourhood for each point and find core points

$k \leftarrow 0$

for $x_i \in Core$ with $id(x_i) = \emptyset$:

    $k \leftarrow k+1$

    $id(x_i) \leftarrow k$

    $DensityConnected(x_i, k)$

> for each unassigned core point, recursively find all of its DR points + assign them all to cluster $k$

$\mathcal{E} \leftarrow \{E_1, \ldots, E_k\}$ where $E_i = \{x \in X \mid id(x) = i\}$

$Noise = \{x \in X \mid id(x) = \emptyset\}$   ← unassigned points are noise points

$\ldots = X \setminus (Core \cup Noise)$ ←

Noise $=$ ...

Border $= X | (\text{Core} \cup \text{Noise})$ ← *points that are not core points or noise points are border points*

return $\Sigma$, Core, Border, Noise

$\text{DensityConnected}(x, k):$ ← *recursively follows neighborhoods of core points to assign to the same cluster*

   for $z \in N_{\varepsilon}(x):$

     $id(z) \leftarrow k$ ← *assign neighborhood points to cluster*

     if $z \in \text{Core}:$

       $\text{DensityConnected}(z, k)$ } *explore the neighborhoods of DC core points to add to the clusters!*

Limitation of DBSCAN: sensitive to $\varepsilon$

   too small → sparse clusters seen as noise

   too large → dense clusters merged

   If there are clusters of different densities, a single $\varepsilon$ may not suffice

      ↘

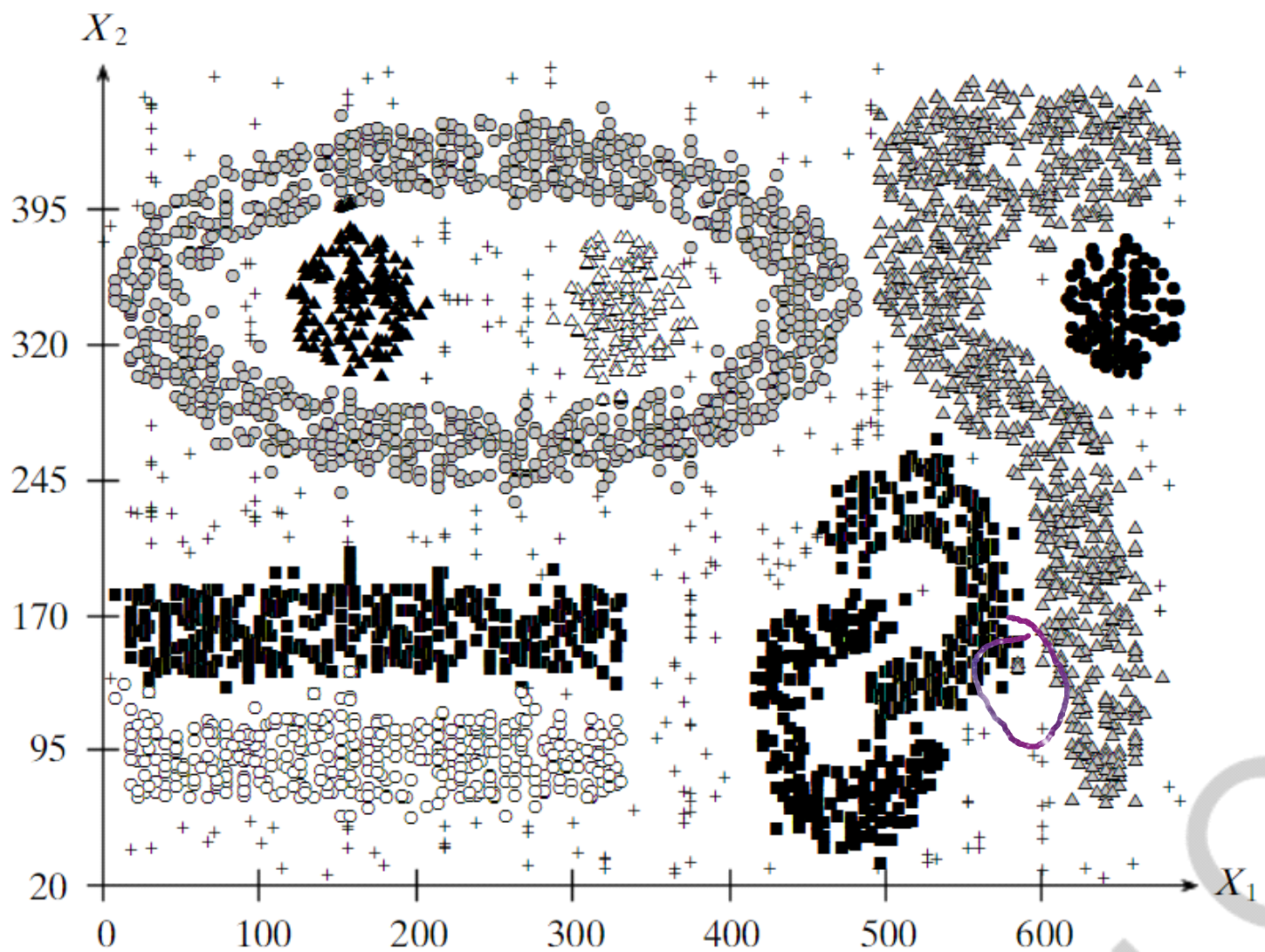      OPTICS uses a <u>range of $\varepsilon$ values</u>

<u>Complexity</u>: $O(n^2)$

Example



**Figure 15.3.** Density-based clusters.