

Hierarchical Clustering Analysis & Density clustering

agglomerative

Analysis of¹ Hierarchical Clustering

Assume that dist comp takes const time

① init pdm $O(n^2)$

② for each merge step 1 less cluster

\Rightarrow at step t recompute $O(n-t)$ dists

\Rightarrow #_{dist} $n-1 \quad n-2 \quad n-3 \quad \dots \quad 1$
step 1 2 3 \dots $n-1$

\Rightarrow total #_{of dists in alg} $(n-1) + (n-2) + (n-3) + \dots + 1 = O(n^2)$

③ at step i we need to:

1. find min dist

2. remove old dist (from individual clusters)

3. add new dist (for the merged cluster)

Let's use Priority Queue implemented w/ a heap (w/ arbitrary deleter)
for heap of size k

1. find min: $O(1)$

2. update: $O(\log k)$

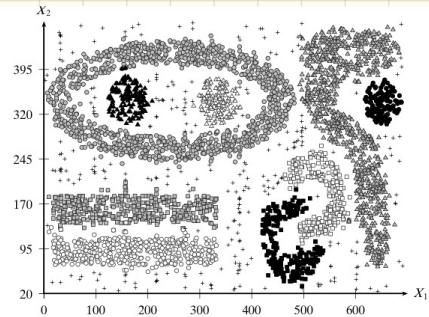
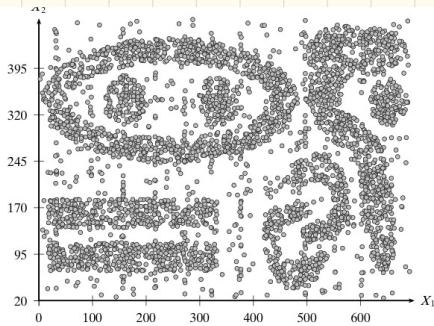
Heap MinCv

A. init the heap w/ $\binom{n}{2}$ elements - $O(n^2)$

B. perform $O(n^2)$ dist updates $\Rightarrow O(n^2 \log n^2) = O(n^2 \log n)$

\Rightarrow total time is $O(n^2 \log n)$

Density based clustering



DBScan

let \$\varepsilon > 0\$ and \$\vec{x} \in \mathbb{R}^d\$

$$\varepsilon\text{-neighborhood of } \vec{x} = B_d(\vec{x}, \varepsilon) = \left\{ \vec{y} \in D \mid \| \vec{y} - \vec{x} \| \leq \varepsilon \right\}$$

using Euclidean

dist but

in general

can use any metric

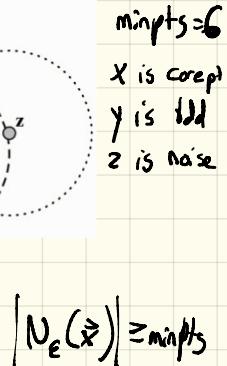


define

$$|N_\varepsilon(\vec{x})| = |\{\vec{y} \in D \mid \vec{y} \in N_\varepsilon(\vec{x})\}|$$

Given param $\text{minpts} \in \mathbb{Z}^+$
for any $\vec{x} \in D$

\vec{x} is a **core point** if there are
at least minpts in its ε -neighborhood



$$|N_\varepsilon(\vec{x})| \geq \text{minpts}$$

\vec{x} is a **border point** if it is
not a core point and \vec{x} is in
 ε -neighborhood of a core point

$$\exists \vec{y} \in D \ni \vec{x} \in N_\varepsilon(\vec{y})$$

and $|N_\varepsilon(\vec{x})| \geq \text{minpts}$

\vec{x} is a **noise point** if not core point nor a border point

Given $\vec{x}, \vec{y} \in D$

- \vec{x} is directly density reachable from \vec{y} if \vec{y} is a core point and $\vec{x} \in N_c(\vec{y})$
- \vec{x} is density reachable from \vec{y} if \exists seq $\vec{x} = \vec{x}_0, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_l = \vec{y}$ w/ \vec{x}_i is directly reachable from \vec{x}_{i-1} for $i=1, \dots, l$
- \vec{x}, \vec{y} density connected if \exists a core point \vec{z} s.t. \vec{y} & \vec{x} are density reachable from \vec{z}
- Density base clustering is find maximal set of density connected pts