

So far: clustering with respect to shape (K-means) and distance between clusters at different scales (hierarchical methods).

A different approach is to consider a notion of "density" as correlating with clusterability.

ex:



The two circles are sample roughly uniformly.
 The density on each is ~~uniform~~ roughly the same, and there is no easy way to stop between them.

"SIGKDD test of time" award



We want a method that can identify these clusters.

DBSCAN: "Density-based spatial clustering of applications with noise"

Let $x \in \mathbb{R}^d$, and let $B_r(x) = \{y \mid \|x-y\|_2 \leq r\}$.

Let $\{x_1, \dots, x_n\} \subset \mathbb{R}^d$ be the data to be clustered. We say x_i is a corepoint if $|B_\epsilon(x)| \geq \text{MinPts}$, where ϵ, MinPts are parameters.

ϵ may be thought of as local neighborhood size: $\epsilon \rightarrow 0 \Rightarrow |B_\epsilon(x)| = 0, \forall x$
 $\epsilon \rightarrow \infty \Rightarrow |B_\epsilon(x)| = n, \forall x$.

MinPts is a threshold for significance: MinPts small suggests every neighborhood is meaningful.

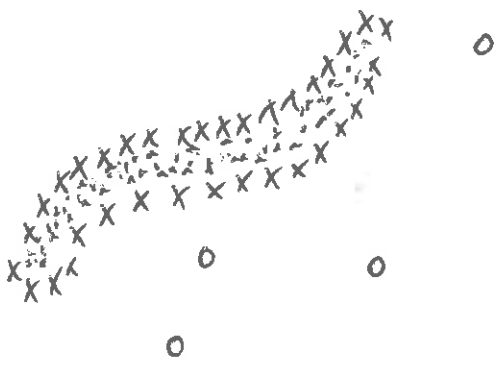
• MinPts large suggests no neighborhood is meaningful. ^②

• A ^{border} ~~core~~ point is x s.t. 1) $|B_\epsilon(x)| \geq \text{MinPts}$

2) $B_\epsilon(z) \ni x$, for some core point z .

• A noise point is every one else.

ex:



\bullet = core point
 x = border point
 o = noise point

} with respect to
some ϵ , MinPts.

• DBSCAN works by linking contiguous regions of "core points" together.

• More precisely, we say:

1) x is directly density reachable from y if a) y is a core point
b) $x \in B_\epsilon(y)$

2) x is density reachable from y if there exists a path

x_0, x_1, \dots, x_L such that: a) $x_0 = x$, $x_L = y$

b.) X_i is directly density reachable from X_{i-1} , $i=2, \dots, L$. ③

3.) X and y are density connected if there exists a core point z such that both X and y are density reachable from z .

• The clusters learned by DBSCAN are the maximal sets of density connected points.

• This method has the benefit of being shape-agnostic, meaning it doesn't favor spherical clusters the way K-means does, for example.

• On the other hand, it depends heavily on ϵ , MinPts, and has large complexity in high dimensions.

• Indeed, computing $B_\epsilon(x)$ requires a nearest neighbors search. A naive implementation is $O(n)$ for each point, so overall $O(n^2)$.

Remark: Indexing structures like k-d trees (~ 1975) and cover trees (~ 2005) speed up their searches to $O(n \log n)$ with a big constant, if certain mathematical assumptions on the data hold.

Note that "density" was implicit here, in the notion of core point, border point, and noise point.

We can also explicitly estimate the density of the points, using a kernel-density estimator.

Let $K: \mathbb{R}^d \rightarrow \mathbb{R}$ be a kernel, i.e. 1) $K(x) \geq 0 \forall x$

2) $\int_{\mathbb{R}^d} K(x) dx = 1$

3) $K(-x) = K(x)$.

Let $p(x) = \frac{1}{n h^d} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$ for some scalar h , which can be

understood as the scale of the density estimator. As $n \rightarrow \infty, \frac{p(x)}{h \rightarrow 0} \int_{\mathbb{R}^d} p(x) dx$

~~estimate~~ estimates the underlying probability distribution of the data (under certain mathematical assumptions we do not discuss).

Choosing $K: \mathbb{R}^d \rightarrow \mathbb{R}$ the kernel is important.

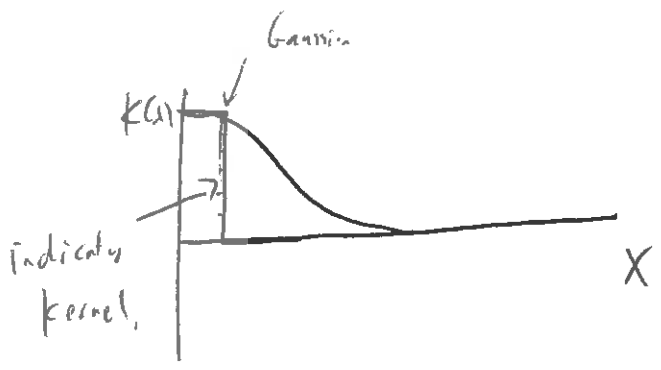
ex. Fix $T > 0$. Let $K_T(x) = \begin{cases} 1, & \|x\|_2 \leq T \\ 0, & \text{else} \end{cases}$
Indicator kernel

• This gives mass wherever $\|x - x_i\| \leq hT$

• To soften the transition, a smooth kernel could be used

ex: $k_{\sigma}(x) = e^{-\|x\|_2^2 / 2\sigma^2}$ for some $\sigma > 0$. This is the Gaussian

kernel.



• Supposing the density is well estimated by $\hat{\rho}(x)$ (happens if n is large, say $n \gg 2^D$ and h small)

This can be used to cluster: associated points to nearby regions of high density.

- Several variations on this theme:
 - DENCLUE
 - Mean-shift clustering
 - Fast search and find of density peaks clustering (FSFOPC)

• Generally good, but may need a lot of points, or computational resources.