

Requirements of Clustering Algorithms

- minimum requirement of domain knowledge
- discovery of clusters with arbitrary shape with good efficiency on large databases

DBSCAN - designed to discover clusters of arbitrary shape relying on density based notion of clusters.

- Requires only one i/p parameter.

Spatial Database is large \Rightarrow automated knowledge discovery.
We solve the classic task of class identification.

Two types of clustering Algorithms - Partitioning and Hierarchical.

Partitioning Algorithms partition the Database D of n objects into k clusters, k being the i/p parameter. Some domain knowledge is required. Usually follows iterative strategy to partition, where each cluster is determined by its gravity centre (KNN) or the point nearest to its center (K-medoid algorithm)

Shape of each cluster in a partitioning algorithm is convex which is very restrictive.

eg. CLARANS

Hierarchical Algorithms create a hierarchical decomposition of D , which is represented by a dendrogram, a tree which keeps splitting the D into smaller subsets until each subset contains only one object, each node of the tree is a cluster.

Leaves to Root \rightarrow (agglomerative approach)

Root to Leaves \rightarrow (divisive approach)

eg. Ecluster, critical distance, $O(n^2)$

A Density Based Notion of Clusters

Database D is k -dimensional space S .

Key-idea - For each point of a cluster the neighborhood of a given radius has to contain a minimum number of points.
Shape of the cluster depends upon the distance funcⁿ chosen, $\text{dist}(P, q)$, between two points p and q .

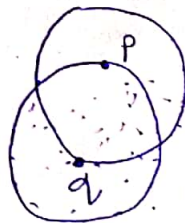
(Eps-neighborhood of a point) - the Eps-neighborhood of a point P denoted by $N_{\text{eps}}(P)$ is defined as

$$N_{\text{eps}}(P) = \{q \in D \mid \text{dist}(p, q) \leq \text{Eps}\}$$

Naive approach would be to set MinPts constant but this method fails \because we have two kinds of points in the cluster core and border, core pts contains more points in their neighborhood than the border points.

- A point p is "directly density-reachable" from a point q wrt $\text{Eps}, \text{MinPts}$ if
 - $p \in N_{\text{eps}}(q)$
 - $|N_{\text{eps}}(q)| \geq \text{MinPts}$ (core point condition)

This condition is ~~sym~~ symmetric iff they belong to core points.



p is border point

q is core point

q is not directly density-reachable from p

- A point p is "density-reachable" from a point q wrt Eps and MinPts if there is a chain of points p_1, \dots, p_n st. $p_1 = q$, $p_n = p$ & p_{i+1} is directly density-reachable from p_i

Travelling does not symmetric, but is symmetric for core points. @AshishSinha

- A point p is density-connected to a point q wrt Eps and $MinPts$ if there is a point O st both p and q are density-reachable from O wrt Eps and $MinPts$.

Symmetric Relation, \Rightarrow Reflexive for density reachable points.

- Cluster is a set of density connected points which is maximal wrt density reachability.

Noise is simply the set of points in D not belonging to any of its clusters.

- Let D be a database of points. A cluster C wrt Eps and $MinPts$ is a non-empty subset of D satisfying the following condition-

- $\forall p, q: \text{if } p \in C \text{ and } q \text{ is density-reachable from } p \text{ wrt } Eps \text{ and } MinPts, \text{ then } q \in C$ (Maximality)
- $\forall p, q \in D: p \text{ is density connected to } q \text{ wrt } Eps \text{ and } MinPts$ (Connectivity)

- Let C_1, \dots, C_k be the clusters of a ~~point~~ Database D , wrt Eps ; and $MinPts_i$ $i=1, \dots, k$, then Noise is the set of points in the database D not belonging to any cluster C_i , i.e. noise = $\{p \in D \mid \forall i: p \notin C_i\}$

Cluster C contains atleast $MinPts$.

Given $Eps, MinPts$ (i) choose arbitrary point in the database following the core point condition i.e. the "seed" (ii) Retrieve all the pts that are directly reachable from ~~pt~~ seed.

- Lemma - if $p \in D$ st $|N_{Eps}(p)| > MinPts$ and $O = \{o \in D \mid \text{st } o \text{ is density-reachable from } p \text{ wrt } Eps \text{ and } MinPts\}$ then O is a cluster wrt Eps and $MinPts$.

Lemma 2: Let C be any cluster in D and let $p \in C$ st $|N_{Eps}(p)| \geq MinPts$. Then C equals to the set $O = \{o \mid o \text{ is density-reachable from } p \text{ w.r.t. } MinPts \text{ and } Eps\}$

i.e. cluster C is uniquely determined by any of its core points

DBSCAN - Density Based Spatial Clustering For Applications with Noise

Uses global values of Eps and $MinPts$ which is determined by the least, dense cluster in the database D .

Algorithm

Starts with an ~~cluster~~ ^{arbitrary} point $p \in D$ and retrieves all the points which are density reachable from p w.r.t Eps and $MinPts$. If p is core pt this procedure produces a cluster & w algorithm moves to other pts $\in D$.

Distance between two sets of pts S_1 and S_2 be defined as $dist(S_1, S_2) = \min\{dist(p, q) \mid p \in S_1, q \in S_2\}$, then the two sets of points having at least the density of the thinnest cluster will be separated from each other iff distance betⁿ them is ~~separated by~~ Eps larger than eps .

DBSCAN (SetOfPoints, Eps , $MinPts$)

// SetOfPoints is unclassified

ClusterID := NextID(Noise);

For i from 1 to SetOfPoints.size Do

Point := SetOfPoints.get(i);

If Point.CID = UNCLASSIFIED THEN

If ExpandCluster(SetOfPoints, Point, ClusterID, Eps , $MinPts$)

ClusterID := NextID(ClusterID)

End If

End If

EndFor

End; // DBSCAN

ExpandCluster (SetOfPoints, Point, CId, Eps, MinPts): Boolean;

Seeds: SetOfPoints.regionQuery (Points, Eps);

If seeds.size < MinPts Then // no core point

SetOfPoint.changeCId (point, NOISE);

Return False;

Else // all points in the seed are density reachable from Point

SetOfPoints.changeCIds (seeds, CId);

Seeds.delete (Point);

While seeds <> Empty Do

currentP := seeds.first();

result := SetOfPoints.regionQuery (currentP, Eps);

If result.size >= MinPts Then

For i from 1 to result.size Do

resultP := result.get(i)

If resultP.CId IN (UNCLASSIFIED, NOISE) Then

If resultP.CId = UNCLASSIFIED Then

Seeds.append (resultP);

EndIf;

SetOfPoints.changeCId (resultP, CId)

EndIf; // Unclassified or Noise.

EndFor;

EndIf; // Result Size >= MinPts

Seeds.delete (current P);

EndWhile; // seed <> Empty

Return True;

EndIf;

End;

SetOfPoints.regionQuery (Points, Eps) - returns eps neighbourhood of points in SetOfPoints as a list of points. (eg use R* trees)

Average runtime complexity of Region Query is $O(\log N)$ and for each of the n points we have at most one region query
 \therefore Average runtime complexity of DBSCAN is $O(N \log N)$

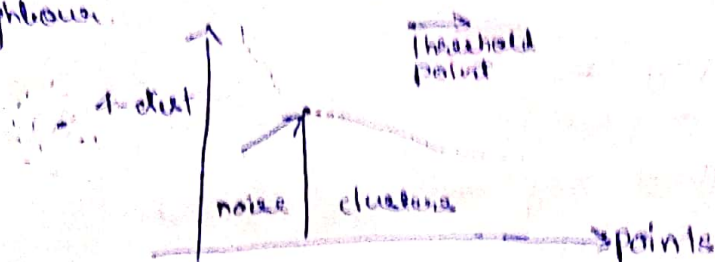
If p is a border point of two clusters C_1 and C_2 , p is assigned to the cluster discovered first.

• Determining the parameters ϵ and MinPts

(of the minimal cluster in the database)

ϵ -distance of a point p to its k^{th} nearest neighbour, then ϵ -neighbourhood of point p contains $k+1$ points for almost all points p .

k -dist: $D \rightarrow \mathbb{R}$ mapping each point to the k^{th} distance from its k^{th} nearest neighbour.



By inspection MinPts parameter is set to 4 - \forall $p \in \text{Database}$.

Steps to computing ϵ and MinPts =

- System computes and displays ϵ -dist graph for the Database
- System derives a proposed threshold pt p by the % of noise entered
- Either accepts or selects another point in the D as the threshold value.

$$\boxed{\epsilon = \epsilon\text{-dist}(p)}$$