

DBSCAN

Paul Antony Job Jacob

9 August 2018

DBSCAN

DBSCAN stands for *Density-Based Spatial Clustering of Applications with Noise*

It is a density-based clustering algorithm which groups together points in a given set that are closely packed together, it also marks outlier's points that lie alone in low-density regions as Noise.

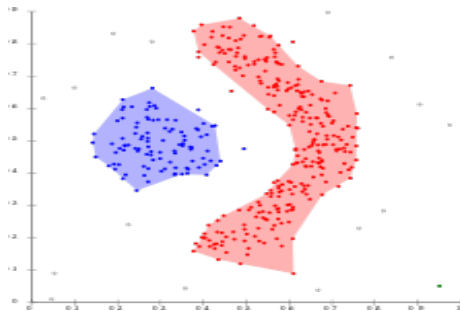


Figure: DBSCAN Clustering

DBSCAN Working

DBSCAN takes two parameters ϵ and MinPtn

ϵ is a distance parameter that defines the radius to search for nearby neighbours and MinPtn minimum number of points in the region required to form a cluster

An n dimensional sphere of radius ϵ is made around a point and points in the sphere are clustered if number of points is greater than MinPtn

There are three types of points in DBSCAN:

- **Core point** – a point that has at least a minimum number of other points (minPts) within its ϵ radius.
- **Border point** – a point is within the ϵ radius of a core point BUT has less than the minimum number of other points (minPts) within its own ϵ radius
- **Noise point** – a point that is neither a core point nor a border point

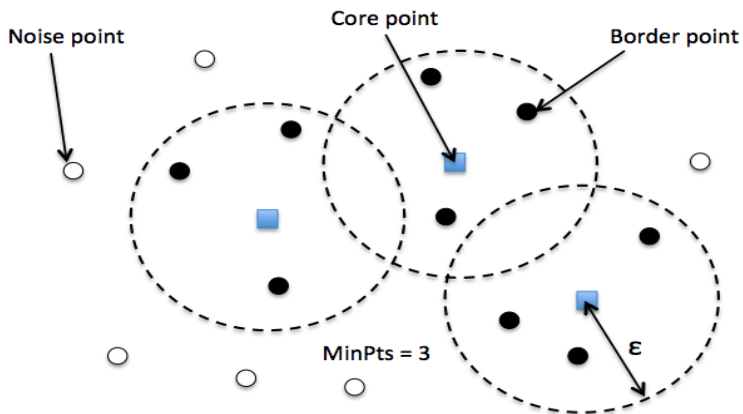


Figure: Cluster Points

Advantages

- DBSCAN does not require one to specify the number of clusters in the data
- It can find arbitrarily shaped cluster. Find cluster completely surrounded by different clusters
- It requires just two parameters and is mostly insensitive to the ordering of the points in the database
- It is robust towards outlier detection (noise)

Disadvantage

- DBSCAN is not entirely deterministic: border points that can be reached by more than one clusters can be part of either cluster depending on the order the data is processed
- If data and scale is unknown choosing a meaningful distance threshold ϵ can be difficult
- Sensitive to clustering parameter MinPoints and ϵ
- DBSCAN cannot cluster data sets well with large differences in densities, since the minPts- ϵ combination cannot then be chosen appropriately for all clusters

Algorithm

INPUT: N objects to be clustered and global parameter ϵ and MinPts

OUTPUT: Clusters of objects

Algorithm

- ① Initialize all object label as *not visited*
- ② Repeat till all points are visited
 - ① Arbitrary select a point P
 - ② Retrieve all neighbour points of P using function `regionQuery(P)`
 - ③ If number of neighbour points is greater than MinPts then P is a core point. Its label is changed to cluster name and cluster is grown using function `growCluster(P)`, else P is labeled as noise

Algorithm growCluster(P)

- ① push P to SearchQueue (SearchQueue is an normal Queue data structure)
- ② while SearchQueue is not empty
 - ① $P = \text{SearchQueue.pop}$
 - ② find neighbour points of P using function regionQuery(P)
 - ③ check each point in neighbour is a core point, change its label to cluster name and push it into SearchQueue.
 - ④ If it's not a noise point change its label to cluster name

Algorithm

Algorithm regionQuery(P)

- ① calculates distance between P and each points P_i in data using function $\text{EuclideanDistance}(P, P_i)$. If distance is less than ϵ then the point P_i is added to list neighbour-points
- ② return neighbour-points

Algorithm EuclideanDistance(P, P_i)

- ① distance between points P and P_i is calculated using the Euclidean Distance equation

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + \dots}$$

Sample Output

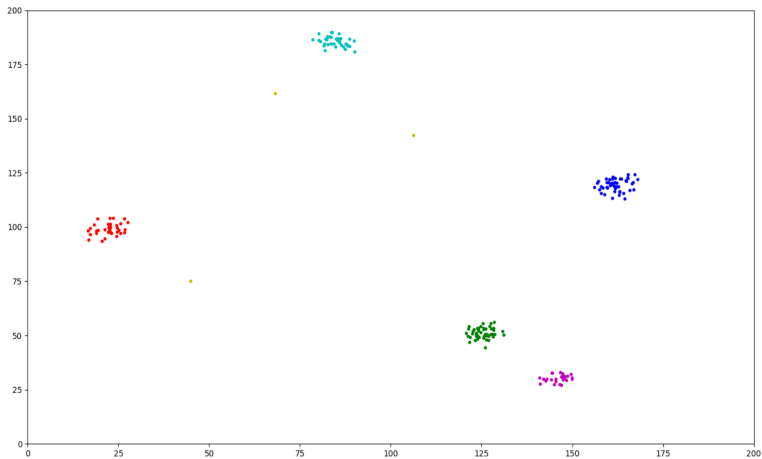


Figure: Sample Output

Reference

- <https://en.wikipedia.org/wiki/DBSCAN>
- <http://mccormickml.com/2016/11/08/dbscan-clustering/>
- <https://towardsdatascience.com/how-dbscan-works-and-why-should-i-use-it-443b4a191c80>
<https://algorithmicthoughts.wordpress.com/2013/05/29/machine-learning-dbscan/>

END