

Open in app



Following ▾

565K Followers



Understanding DBSCAN Algorithm and Implementation from Scratch

DBSCAN Algorithm Step by Step, Python Implementation, and Visualization.



Andrewngai Jun 9, 2020 · 5 min read ★

What is DBSCAN

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a commonly used unsupervised clustering algorithm proposed in 1996. Unlike the most well known K-mean, DBSCAN does not need to specify the number of clusters. It can automatically detect the number of clusters based on your input data and parameters. More importantly, DBSCAN can find arbitrary shape clusters that k-means are not able to find. For example, a cluster surrounded by a different cluster.

Click to add this story to list.

Got it



DBSCAN vs K-means, [credit](#)

Also, DBSCAN can handle noise and outliers. All the outliers will be identified and marked without being classified into any cluster. Therefore, DBSCAN can also be used for Anomaly Detection (Outlier Detection)

Before we take a look at the pseudocode, we need to first understand some basic concepts and terms. Eps, Minpts, Directly density-reachable, density-reachable, density-connected, core point and border point

First of all, there are two parameters we need to set for DBSCAN, Eps, and MinPts.

Eps: Maximum radius of the neighborhood

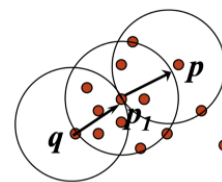
MinPts: Minimum number of points in an Eps-neighbourhood of that point

And there is the concept of **Directly density-reachable**: A point p is directly density reachable from a point q w.r.t. Eps, MinPts, if $N_{Eps}(q) = \{p \text{ belongs to } D \mid \text{dist}(p, q) \leq Eps\}$ and $|N_{Eps}(q)| \geq \text{MinPts}$. Let's take a look at an example with $\text{Minpts} = 5$, $Eps =$



■ Density-reachable:

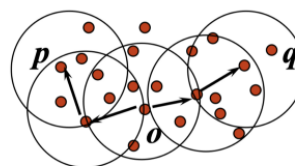
- A point p is **density-reachable** from a point q w.r.t. Eps , $MinPts$ if there is a chain of points p_1, \dots, p_n , $p_1 = q$, $p_n = p$ such that p_{i+1} is directly density-reachable from p_i



Density-reachable example

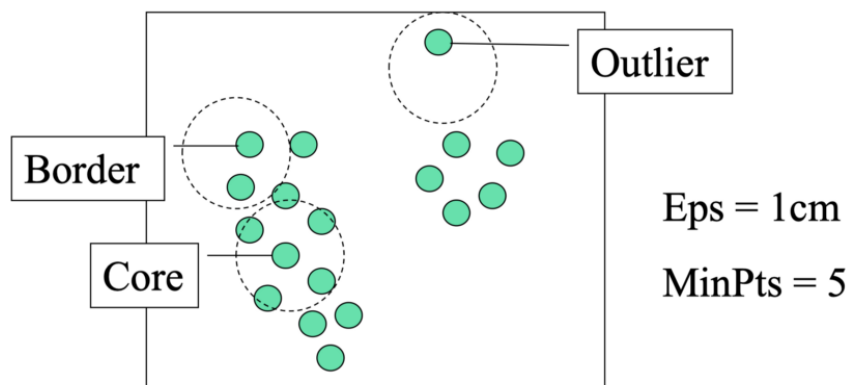
■ Density-connected

- A point p is **density-connected** to a point q w.r.t. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o w.r.t. Eps and $MinPts$



Density-connected example

Finally, a point is a core point if it has more than a specified number of points ($MinPts$) within Eps . These are points that are at the interior of a cluster A. And a border point has fewer than $MinPts$ within Eps , but is in the neighborhood of a core point. We can also define the outlier(noise) point, which is the points that are neither core nor border points.



Core point, Border point, Outlier Point examples

Now, let's take a look at how DBSCAN algorithm actually works. Here is the pseudocode.

1. Arbitrary select a point p
2. Retrieve all points density-reachable from p based on Eps and $MinPts$
3. If p is a core point, a cluster is formed
4. If p is a border point, no points are density-reachable from p and DBSCAN visits the next point of the database

[Open in app](#)

If a spatial index is used, the computational complexity of DBSCAN is $O(n \log n)$, where n is the number of database objects. Otherwise, the complexity is $O(n^2)$

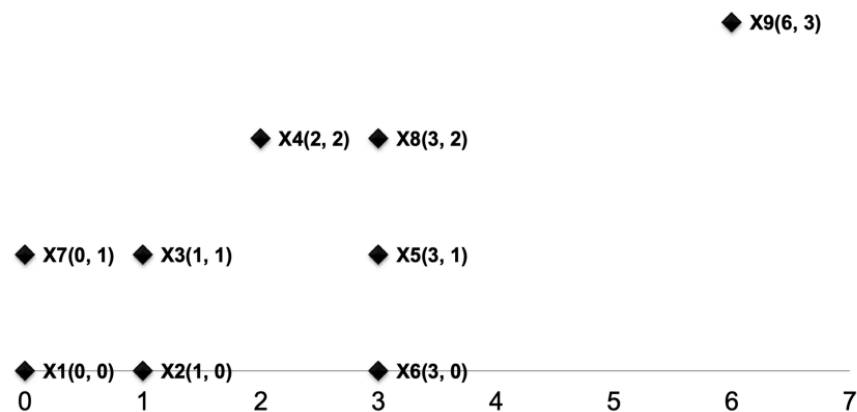
Example

Consider the following 9 two-dimensional data points:

$x_1(0,0)$, $x_2(1,0)$, $x_3(1,1)$, $x_4(2,2)$, $x_5(3,1)$, $x_6(3,0)$, $x_7(0,1)$, $x_8(3,2)$, $x_9(6,3)$

Use the Euclidean Distance with $Eps = 1$ and $MinPts = 3$. Find all core points, border points and noise points, and show the final clusters using DBSCAN algorithm. Let's show the result step by step.

Data Points



Example Data Visualization

First, Calculate the $N(p)$, Eps-neighborhood of point p

$$N(x_1) = \{x_1, x_2, x_7\}$$

$$N(x_2) = \{x_2, x_1, x_3\}$$

$$N(x_3) = \{x_3, x_2, x_7\}$$

$$N(x_4) = \{x_4, x_8\}$$

$$N(x_5) = \{x_5, x_6, x_8\}$$

$$N(x_6) = \{x_6, x_5\}$$

$$N(x_7) = \{x_7, x_1, x_3\}$$

$$N(x_8) = \{x_8, x_4, x_5\}$$

$$N(x_9) = \{x_9\}$$

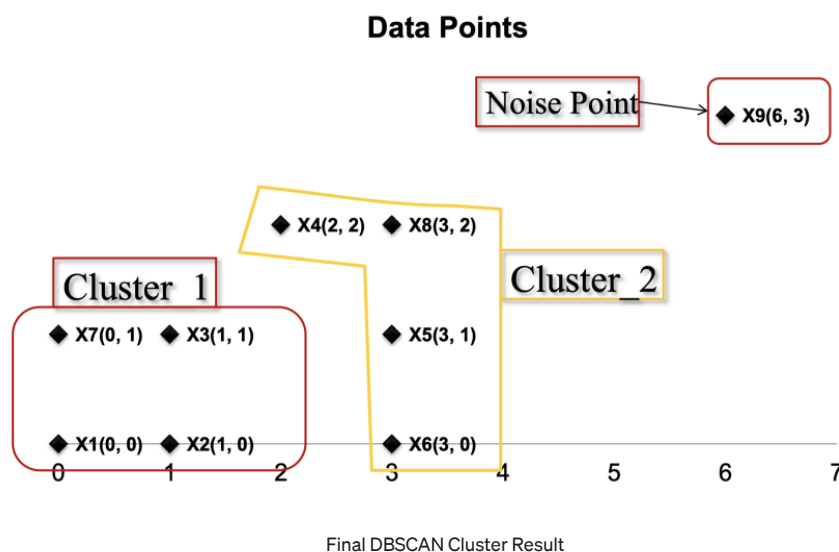
If the size of $N(p)$ is at least $MinPts$, then p is said to be a core point. Here the given $MinPts$ is 3, thus the size of $N(p)$ is at least 3. **Thus core points are: $\{x_1, x_2, x_3, x_5, x_7, x_8\}$**



$N(x_6) = \{x_6, x_5\}$. here x_8 and x_5 are core points, So both x_4 and x_6 are border points. Obviously, the point left, x_9 is a noise point.

Now, let's follow the predecode to produce the clusters.

1. Arbitrary select a point p , now we choose x_1
2. Retrieve all points density-reachable from x_1 : $\{x_2, x_3, x_7\}$
3. Here x_1 is a core point, a cluster is formed. So we have Cluster_1: $\{x_1, x_2, x_3, x_7\}$
4. Next, we choose x_5 , Retrieve all points density-reachable from x_5 : $\{x_8, x_4, x_6\}$
5. Here x_5 is a core point, a cluster is formed. So we have Cluster_2: $\{x_5, x_4, x_8, x_6\}$
6. Next, we choose x_9 , x_9 is a noise point, noise points do NOT belong to any clusters.
7. Thus the algorithm stops here.



Python Implementation

Here is some sample code to build FP-tree from scratch and find all frequency itemsets in Python 3. I have also added visualization of the points and marked all outliers in blue.

```

1  import numpy as np
2  import collections
3  import matplotlib.pyplot as plt
4  import queue
5  import scipy.io as spio
6
7  #Define label for differnt point group
8  NOISE = 0
9  UNASSIGNED = 0
10 core=-1
11 edge=-2
12
13
14
15 #function to find all neighbor points in radius
16 def neighbor_points(data, pointId, radius):
17     points = []
18     for i in range(len(data)):

```

Open in app



```

22     return points
23
24 #DB Scan algorithm
25 def dbscan(data, Eps, MinPt):
26     #initilize all pointlable to unassign
27     pointlabel = [UNASSIGNED] * len(data)
28     pointcount = []
29     #initilize list for core/noncore point
30     corepoint=[]
31     noncore=[]
32
33     #Find all neighbor for all point
34     for i in range(len(data)):
35         pointcount.append(neighbor_points(train,i,Eps))
36
37     #Find all core point, edgepoint and noise
38     for i in range(len(pointcount)):
39         if (len(pointcount[i])>=MinPt):
40             pointlabel[i]=core
41             corepoint.append(i)
42         else:
43             noncore.append(i)
44
45     for i in noncore:
46         for j in pointcount[i]:
47             if j in corepoint:
48                 pointlabel[i]=edge
49
50         break
51
52     #start assigning point to luster
53     cl = 1
54     #Using a Queue to put all neighbor core point in queue and find neighbor's neighbor
55     for i in range(len(pointlabel)):
56         q = queue.Queue()
57         if (pointlabel[i] == core):
58             pointlabel[i] = cl
59             for x in pointcount[i]:
60                 if(pointlabel[x]==core):
61                     q.put(x)
62                     pointlabel[x]=cl
63                 elif(pointlabel[x]==edge):
64                     pointlabel[x]=cl
65             #Stop when all point in Queue has been checked
66             while not q.empty():
67                 neighbors = pointcount[q.get()]
68                 for y in neighbors:
69                     if (pointlabel[y]==core):
70                         pointlabel[y]=cl
71                         q.put(y)
72                     if (pointlabel[y]==edge):
73                         pointlabel[y]=cl
74             cl=cl+1 #move to next cluster
75
76     return pointlabel,cl
77
78 #Function to plot final result
79 def plotRes(data, clusterRes, clusterNum):
80     nPoints = len(data)
81     scatterColors = ['black', 'green', 'brown', 'red', 'purple', 'orange', 'yellow']
82     for i in range(clusterNum):
83         if (i==0):
84             #Plot all noise point as blue
85             color='blue'

```

[Open in app](#)

```
89     for j in range(nPoints):
90         if clusterRes[j] == i:
91             x1.append(data[j, 0])
92             y1.append(data[j, 1])
93     plt.scatter(x1, y1, c=color, alpha=1, marker='.')
94
95
96 #Load Data
97 raw = spio.loadmat('DBSCAN.mat')
98 train = raw['Points']
99
100 #Set EPS and Minpoint
101 epss = [5,10]
102 minptss = [5,10]
103 # Find ALL cluster, outliers in different setting and print resultsw
104 for eps in epss:
105     for minpts in minptss:
106         print('Set eps= ' + str(eps)+ ', Minpoints = '+str(minpts))
107         pointlabel,cl = dbscan(train,eps,minpts)
108         plotRes(train, pointlabel, cl)
109         plt.show()
110         print('number of cluster found: ' + str(cl-1))
111         counter=collections.Counter(pointlabel)
112         print(counter)
113         outliers = pointlabel.count(0)
114         print('numbrer of outliers found: '+str(outliers) +'\n')
```

DBSCAN.py hosted with ❤ by GitHub

[view raw](#)

Thanks for reading and I am looking forward to hearing your questions and thoughts. If you want to learn more about Data Science and Cloud Computing, you can find me on [Linkedin](#).

Photo by [Alfons Morales](#) on [Unsplash](#)*Reference*<https://github.com/NSHipster/DBSCAN>

[Open in app](#)

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

Get this newsletter

Emails will be sent to xili9@mail.yu.edu.
[Not you?](#)

[Dbscan](#) [Data Mining](#) [Clustering](#) [Python](#) [Outlier Detection](#)

[About](#) [Write](#) [Help](#) [Legal](#)

Get the Medium app

