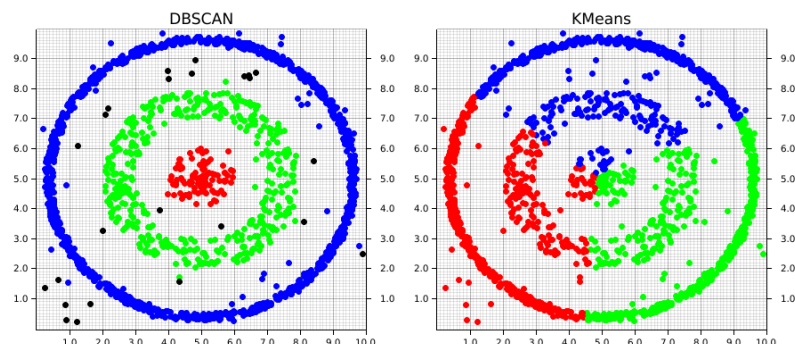




## 1. Introduction to DBSCAN

**DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** is a density-based clustering algorithm used to identify clusters in datasets based on the density of data points. It is particularly useful for discovering clusters of arbitrary shape and handling noise (outliers).

Unlike K-Means, DBSCAN does not require specifying the number of clusters in advance and can identify outliers as points that do not belong to any cluster.



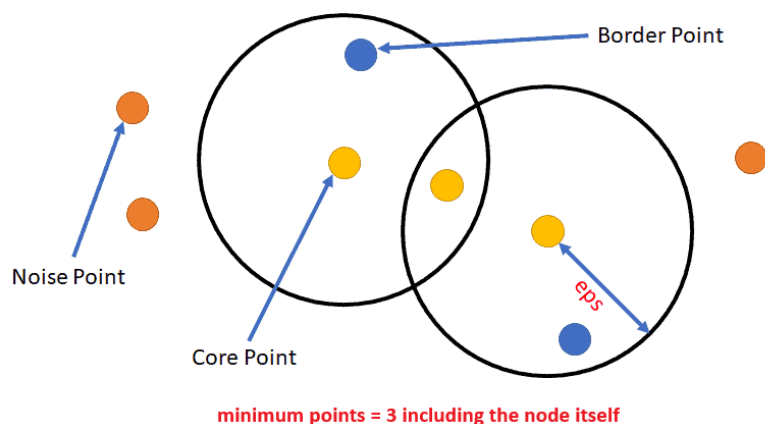
## 2. How Does DBSCAN Work

DBSCAN groups points that are closely packed together, marking as outliers points that lie alone in low-density regions. It works based on two key parameters:

- **Epsilon ( $\epsilon$ ):** The maximum distance between two points for them to be considered as part of the same neighborhood.
- **MinPoints:** The minimum number of points required to form a dense region (i.e., a cluster).

Key concepts in DBSCAN:

- **Core Point:** A point with at least MinPoints (**Including Itself**) within a radius  $\epsilon$ .
- **Border Point:** A point that has fewer than MinPoints within  $\epsilon$ , but is in the neighborhood of a core point.
- **Noise Point (Outlier):** A point that is neither a core nor a border point.



**DBSCAN Steps:**

- For each point in the dataset, identify its neighborhood based on  $\epsilon$  and MinPoints.
- If a point is a core point, form a cluster around it and recursively expand the cluster by including all density-reachable points.
- Continue until all points have been visited, and label remaining points that do not belong to any cluster as noise.

DBSCAN creates clusters by iteratively adding points that are density-reachable from core points.

### The Concept of Density-Reachable (directly and indirectly):

A point  $y$  is said to be *density-reachable* from  $x$  if there is a path  $(p_1, \dots, p_n)$  with  $(p_1 = x)$  and  $(p_n = y)$ , such that each point on the path (except possibly  $(p_n)$ ) is a **core** point and each step  $(p_i \rightarrow p_{i+1})$  is *directly density-reachable*.

An object  $y$  is **directly density-reachable** from object  $x$  if  $x$  is a core object and  $y$  lies in  $x$ 's epsilon ( $\epsilon$ )-neighborhood.

In the image (MinPts = 3, including the point itself), **b**, **c**, and **d** are core points, **a** and **e** are border points, and **f**, **g**, **h** are noise points. From this:

- **a** is directly density-reachable from **b**.
- **b** is directly density-reachable from **c**.
- Hence **a** is (indirectly) density-reachable from **d** via the chain  $(d \rightarrow c \rightarrow b \rightarrow a)$ .
- **d** is not density-reachable from **a**, since **a** is not a core point, so density-reachability is not symmetric.

DBSCAN in Action  YouTube

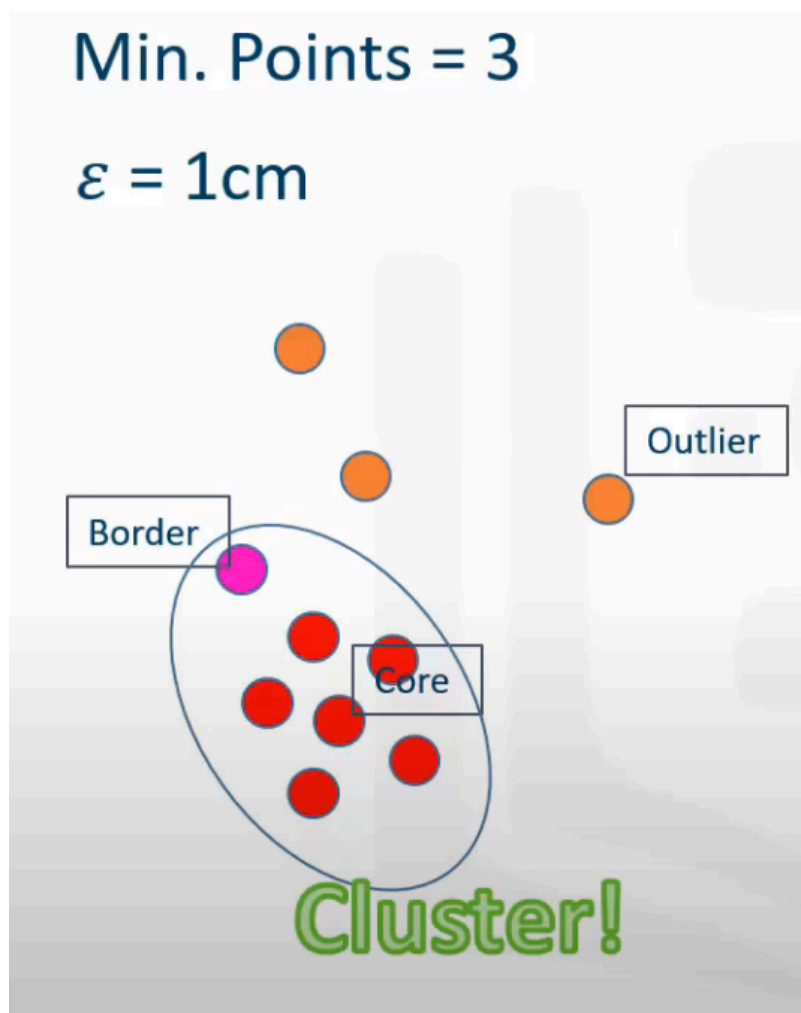
### YouTube Example: How Does the Algorithm Work?

Min. Points = 3  $\epsilon$  = 1 cm

Step 1: Pick a random point that has not yet been assigned to a cluster or designated as an Outlier. Determine if it is a Core Point. If not, label the point as Outlier.

Step 2: Once a Core Point has been found, add all directly reachable to its cluster. Then do neighbor jumps to each reachable point and add them to the cluster. If an Outlier has been added, label it as a border point.

Step 3: Repeat these steps until all points are assigned a cluster or label as Outlier.





- ✓ A Detailed example showing the calculations performed by the algorithm:

We will work with:

- **Epsilon ( $\epsilon$ ): 1.9**
- **MinPoints: 4**

#### i. Compute the Distance for Each Point

- In this step, we compute the Euclidean distances between all points in the dataset, creating a **lower triangular distance matrix**.
- Each cell in the matrix represents the distance between two points, while the diagonal is excluded since it represents the distance of a point to itself. This matrix forms the basis for identifying the neighbors of each point within the epsilon ( $\epsilon = 1.9$ ) radius.

Points	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12
P1: (3,7)												
P2: (4,6)	1.41											
P3: (5,5)	2.83	1.41										
P4: (6,4)	4.24	2.83	1.41									
P5: (7,3)	5.66	4.24	2.83	1.41								
P6: (6,2)	5.83	4.47	3.16	2.24	1.00							
P7: (7,2)	6.40	5.00	3.61	2.24	1.00	1.00						
P8: (8,4)	5.83	4.47	3.16	2.24	1.41	2.24	1.41					
P9: (3,3)	5.00	4.24	3.16	3.16	4.12	3.00	3.16	4.47				
P10: (2,6)	1.41	2.00	3.16	4.24	5.39	5.83	6.40	6.32	3.00			
P11: (3,5)	2.00	1.41	2.00	3.16	4.24	4.47	5.00	4.47	2.24	1.00		
P12: (2,4)	3.16	2.83	3.16	4.00	5.10	5.39	5.83	5.10	2.00	2.00	1.41	0

#### How to Read the Triangular Distances Matrix

Points	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12
P1: (3,7)												
P2: (4,6)	1.41											
P3: (5,5)	2.83	1.41										
P4: (6,4)	4.24	2.83	1.41									
P5: (7,3)	5.66	4.24	2.83	1.41								
P6: (6,2)	5.83	4.47	3.16	2.24	1.00							
P7: (7,2)	6.40	5.00	3.61	2.24	1.00	1.00						
P8: (8,4)	5.83	4.47	3.16	2.24	1.41	2.24	1.41					
P9: (3,3)	5.00	4.24	3.16	3.16	4.12	3.00	3.16	4.47				
P10: (2,6)	1.41	2.00	3.16	4.24	5.39	5.83	6.40	6.32	3.00			
P11: (3,5)	2.00	1.41	2.00	3.16	4.24	4.47	5.00	4.47	2.24	1.00		
P12: (2,4)	3.16	2.83	3.16	4.00	5.10	5.39	5.83	5.10	2.00	2.00	1.41	0

#### ii. Count the Number of Neighbors for Each Point

- After computing the distances for each point, we list the neighbors of each point within the  $\epsilon = 1.9$  radius and counting them.
- This information will be used to classify points as following:
  - a core point requires at least 4 neighbors

- a border point is connected to a core point but has fewer than 4 neighbors
- noise points which are the nodes that has not no neighbors
- For example:
  - P1 has neighbors P2 and P10, as seen in the the table below.
  - P5 has 4 neighbors (P4, P6, P7, P8), while P8 has only 1 neighbor (P5) and P9 has no neighbors.

Point	Neighbors ( $\epsilon=1.9$ )	# Neighbors
P1: (3,7)	P2, P10	2
P2: (4,6)	P1, P3, P11	3
P3: (5,5)	P2, P4	2
P4: (6,4)	P3, P5	2
P5: (7,3)	P4, P6, P7, P8	4
P6: (6,2)	P5, P7	2
P7: (7,2)	P5, P6	2
P8: (8,4)	P5	1
P9: (3,3)	-	0
P10: (2,6)	P1, P11	2
P11: (3,5)	P2, P10, P12	3
P12: (2,4)	P9, P11	2

### iii. Classify Points into Core, Border or Noise

In this step, we classify each point in the dataset as either **Core, Border, or Noise based on the number of neighbors within the epsilon ( $\epsilon = 1.9$ ) radius:**

a. **Core Points:** Points that have at least 4 neighbors within the  $\epsilon$  radius (including the point itself). For example:

- **P2** is a core point because it has 3 neighbors: P1, P3, and P11 which makes them 4 points including the point itself.
- **P5** is a core point because it has 4 neighbors: P4, P6, P7, and P8 which makes them 5 points including the point itself.
- **P11** is a core point because it has 3 neighbors (P2, P10, and P12) and is connected to the core point P2 which makes them 4 points including the point itself.

b. **Border Points:** Points that have fewer than 4 neighbors but are within the  $\epsilon$  radius of at least one core point. For example:

- **P1** is a border point because it has only 2 neighbors (P2 and P10) but is connected to the core point P2.

c. **Noise Points:** Points that have no neighbors within the  $\epsilon$  radius and are not connected to any core point. For example:

- **P9** is classified as noise because it has no neighbors.

This step ensures that each point is categorized appropriately, which is essential for defining the clusters in the DBSCAN algorithm. Core points form the dense regions of clusters, border points connect to these dense regions, and noise points are outliers.

Point	Neighbors ( $\epsilon=1.9$ )	Status
P1: (3,7)	P2, P10	Border
P2: (4,6)	P1, P3, P11	<b>Core</b>
P3: (5,5)	P2, P4	Border
P4: (6,4)	P3, P5	Border
P5: (7,3)	P4, P6, P7, P8	<b>Core</b>
P6: (6,2)	P5, P7	Border
P7: (7,2)	P5, P6	Border
P8: (8,4)	P5	Border
P9: (3,3)	-	Noise
P10: (2,6)	P1, P2, P11	Border
P11: (3,5)	P2, P10, P12	<b>Core</b>
P12: (2,4)	P9, P11	Border

As demonstrated by the table above, the DBSCAN algorithm identified **three core nodes** and **one noise point** based on the parameters

$\epsilon = 1.9$  and  $\text{minPts} = 4$ ,

Nevertheless, DBSCAN will result in **TWO** clusters as **P2 and P11** will be joined into a single cluster.

- **Cluster 0:** Includes points P1, P2, P3, P10, P11, and P12. These points form a dense region in the upper-left area of the diagram.
- **Cluster 1:** Includes points P4, P5, P6, P7, and P8. This cluster represents another dense region in the lower-right part of the diagram.
- **Noise Point:** P9 is classified as noise because it does not have any neighbors within the  $\epsilon$  radius and is not connected to any core point. It is treated as an outlier.

### 3. DBSCAN in Python

#### Python - Example 1

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import DBSCAN
5
6 # Step 1: Define the dataset
7 points = np.array([
8     [3, 7], # P1
9     [4, 6], # P2
10    [5, 5], # P3
11    [6, 4], # P4
12    [7, 3], # P5
13    [6, 2], # P6
14    [7, 2], # P7
15    [8, 4], # P8
16    [3, 3], # P9
17    [2, 6], # P10
18    [3, 5], # P11
19    [2, 4], # P12
20 ])
21
22 # Step 2: Set DBSCAN parameters
23 epsilon = 1.9
24 min_pts = 4
25
26 # Step 3: Apply DBSCAN
27 dbscan = DBSCAN(eps=epsilon, min_samples=min_pts)
28 labels = dbscan.fit_predict(points)
29
30 # Create a DataFrame to hold the points and their cluster labels
31 df = pd.DataFrame(points, columns=['X1', 'X2'])
32 df['Cluster'] = labels # Add cluster labels to the DataFrame
33 df

```

	X1	X2	Cluster
0	3	7	0
1	4	6	0
2	5	5	0
3	6	4	1
4	7	3	1
5	6	2	1
6	7	2	1
7	8	4	1
8	3	3	-1
9	2	6	0
10	3	5	0
11	2	4	0

Next steps: [Generate code with df](#) [New interactive sheet](#)

#### Print out core points

```

1 # Print out the core points
2 core_points = points[dbscan.core_sample_indices_]
3 print("Core points:\n", core_points)

```

```

Core points:
[[4 6]
 [7 3]]

```

[3 5]]




**Print Out The Classification of All Nodes, Core, Border or Noise**

The code below is for demonstration only i.e. studying this code below is optional

```

1 core_samples_mask = np.zeros_like(labels, dtype=bool)
2 core_samples_mask[dbscan.core_sample_indices_] = True
3
4 # Assign point types
5 point_types = []
6 for i in range(len(labels)):
7     if labels[i] == -1:
8         point_types.append('Noise')
9     elif core_samples_mask[i]:
10        point_types.append('Core')
11    else:
12        point_types.append('Border')
13
14 # Add to DataFrame
15 df['Point_Type'] = point_types
16 df

```

	X1	X2	Cluster	Point_Type	
0	3	7	0	Border	
1	4	6	0	Core	
2	5	5	0	Border	
3	6	4	1	Border	
4	7	3	1	Core	
5	6	2	1	Border	
6	7	2	1	Border	
7	8	4	1	Border	
8	3	3	-1	Noise	
9	2	6	0	Border	
10	3	5	0	Core	
11	2	4	0	Border	

Next steps: [Generate code with df](#) [New interactive sheet](#)

**Validate the Triangular Distances Matrix**

The code below is for demonstration only i.e. studying this code below is optional

```

1
2 from sklearn.neighbors import NearestNeighbors
3
4 # Use NearestNeighbors to find neighbors within the epsilon radius
5 nbrs = NearestNeighbors(radius=epsilon)
6 nbrs.fit(points)
7 all_neighbors = nbrs.radius_neighbors(points, return_distance=False)
8
9 # Build the table of neighbors
10 point_names = [f'P{i+1}': ({x[0]}, {x[1]}) for i, x in enumerate(points)]
11 neighbors_list = []
12 neighbor_counts = []
13
14 for i, neighbors in enumerate(all_neighbors):
15     # Remove self from neighbor list
16     neighbor_indices = neighbors[neighbors != i]
17     neighbor_labels = [f'P{j+1}' for j in neighbor_indices]
18     neighbors_list.append(' '.join(neighbor_labels))
19     neighbor_counts.append(len(neighbor_indices))
20
21 # Create DataFrame

```

```

22 neighbors_df = pd.DataFrame({
23     'Point': point_names,
24     'Neighbors ( $\epsilon=1.9$ )': neighbors_list,
25     '# Neighbors': neighbor_counts
26 })
27
28 neighbors_df
29

```

	Point	Neighbors ( $\epsilon=1.9$ )	# Neighbors
0	P1: (3,7)	P2, P10	2
1	P2: (4,6)	P1, P3, P11	3
2	P3: (5,5)	P2, P4	2
3	P4: (6,4)	P3, P5	2
4	P5: (7,3)	P4, P6, P7, P8	4
5	P6: (6,2)	P5, P7	2
6	P7: (7,2)	P5, P6	2
7	P8: (8,4)	P5	1
8	P9: (3,3)	P12	1
9	P10: (2,6)	P1, P11	2
10	P11: (3,5)	P2, P10, P12	3
11	P12: (2,4)	P9, P11	2

Next steps: [Generate code with neighbors\\_df](#) [New interactive sheet](#)

#### Plot the Clusters

```

1  # Plot the clusters
2  plt.figure(figsize=(8, 6))
3  plt.scatter(df['X1'], df['X2'], c=df['Cluster'], cmap='rainbow', marker='o', s=100)
4  plt.title(f'DBSCAN Clustering ( $\epsilon=\text{epsilon}$ , minPts={min_pts})')
5  plt.xlabel('X-axis')
6  plt.ylabel('Y-axis')
7
8  # Annotate points
9  for i, (x, y) in enumerate(zip(df['X1'], df['X2'])):
10     plt.text(x, y, f'P{i+1}', fontsize=10, ha='right')
11
12  plt.grid(True)
13  plt.show()
14

```

DBSCAN Clustering ( $\epsilon=1.9$ , minPts=4)

Explanation of the Code:

1. **Dataset Creation:** The `points` array defines the dataset with coordinates for each point (e.g., P1 = (3,7), P2 = (4,6)).

2. **DBSCAN Implementation:**

- The `DBSCAN` class from `sklearn.cluster` is used with `eps=1.9` ( $\epsilon$ ) and `min_samples=4` (minPts).
- The `fit_predict()` method clusters the data and assigns labels to each point.

3. **Printing Results:**

- Each point's cluster label is printed, with noise points labeled as `Noise`.

4. **Visualization:**

- The scatter plot shows the clusters with distinct colors and noise points in black. Annotations (P1, P2, etc.) make it easy to identify each point.

4. **Advantages and Disadvantages of DBSCAN**

**Advantages:**

- Does not require specifying the number of clusters in advance.
- Can find arbitrarily shaped clusters.
- Identifies outliers as noise points.

**Disadvantages:**

- Performance depends on choosing appropriate values for  $\epsilon$  and MinPoints.
- Struggles with varying densities of clusters.
- High-dimensional data can make it difficult to choose  $\epsilon$ .

