

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

I302 - Aprendizaje Automático y Aprendizaje Profundo

DBSCAN es un algoritmo de clustering que requiere dos parámetros clave: ϵ (el radio de la vecindad) y *minPts* o K (el número mínimo de puntos necesarios para que una región se considere densa). El algoritmo funciona de la siguiente manera:

1. **Seleccionar un punto arbitrario:** El algoritmo comienza seleccionando un punto arbitrario p en el conjunto de datos.
2. **Examinar la ϵ -vecindad de p :** Se visita la vecindad de radio ϵ alrededor del punto p . Si esta vecindad contiene al menos *minPts* puntos (incluyendo al propio punto p), se considera que p está en una región densa y se inicia un nuevo clúster C . Si no hay suficientes puntos en la vecindad, p se etiqueta como ruido. Es importante notar que un punto etiquetado inicialmente como ruido puede más tarde ser parte de un clúster si se encuentra en la vecindad de otro punto denso.
3. **Expandir el clúster:** Si p es parte de un clúster, todos los puntos en su ϵ -vecindad también se consideran parte del clúster. A continuación, se examina la vecindad de cada uno de estos puntos. Si alguno de estos puntos tiene al menos *minPts* puntos en su ϵ -vecindad, estos puntos adicionales se añaden al clúster y su vecindad también se explora. Este proceso se repite de manera iterativa, expandiendo el clúster hasta que no se puedan añadir más puntos.
4. **Repetir el proceso:** Una vez que el clúster actual se ha expandido por completo, se selecciona un nuevo punto no visitado y se repite el proceso para descubrir nuevos clústeres o identificar más puntos de ruido.

Se recomienda consultar el paper de DBSCAN para una explicación más detallada del algoritmo y sus propiedades [1]. El pseudocódigo del algoritmo DBSCAN se muestra en el **Algoritmo 1**.

Algoritmo 1 DBSCAN

```
1: Input: Dataset  $D$ , distance function  $dist$ , minimum points  $minPts$ ,  
   radius  $\epsilon$   
2: Output: Clusters  $C$   
3:  $C \leftarrow \{\}$  ▷ Initialize an empty set of clusters  
4:  $visited \leftarrow \{\}$  ▷ Initialize an empty set of visited points  
5: for each point  $p \in D$  do  
6:   if  $p \in visited$  then  
7:     continue ▷ Skip  $p$  if it has already been visited  
8:   end if  
9:    $visited \leftarrow visited \cup \{p\}$  ▷ Mark  $p$  as visited  
10:   $neighbors \leftarrow \text{REGIONQUERY}(p, \epsilon)$   
11:  if  $|neighbors| < minPts$  then  
12:    mark  $p$  as noise  
13:  else  
14:     $C \leftarrow C \cup \{\text{EXPANDCLUSTER}(p, neighbors, C)\}$   
15:  end if  
16: end for  
  
17: function  $\text{EXPANDCLUSTER}(p, neighbors, C)$   
18:    $cluster \leftarrow \{p\}$  ▷ Initialize the new cluster with  $p$   
19:   for each point  $n \in neighbors$  do  
20:     if  $n \notin visited$  then  
21:        $visited \leftarrow visited \cup \{n\}$  ▷ Mark  $p$  as visited  
22:        $newNeighbors \leftarrow \text{REGIONQUERY}(n, \epsilon)$   
23:       if  $|newNeighbors| \geq minPts$  then  
24:          $neighbors \leftarrow neighbors \cup newNeighbors$  ▷ Add neighbors  
25:       end if  
26:     end if  
27:     if  $n \notin anyCluster$  then ▷ If neighbor is not part of any cluster  
28:        $cluster \leftarrow cluster \cup \{n\}$  ▷ Add it to the current cluster  
29:     end if  
30:   end for  
31:   return  $cluster$   
32: end function  
  
33: function  $\text{REGIONQUERY}(p, \epsilon)$   
34:   return  $\{q \in D \mid dist(p, q) \leq \epsilon\}$   
35: end function
```

Referencias

- [1] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, 1996, pp. 226-231.