

## DEFINIZIONE DI DENSITA'

Per poter definire cluster basati sulla densità bisogna definire come si intende densità:

Abbiamo due parametri:

- 1)  $\epsilon$ : Massimo raggio per cui definire il vicinato
- 2)  $\text{MinPoints}$ : Minimo numero di punti per far sì che il vicinato sia denso.

### CORE OBJECT

È un oggetto che entro il raggio  $\epsilon$  ha un vicinato che comprende almeno  $\text{MinPoints}$ .

### DIRECTLY DENSITY REACHABLE

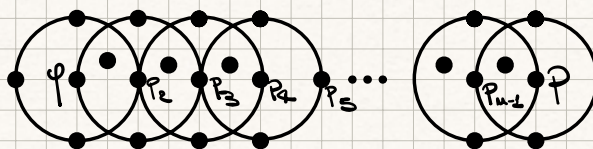
Un oggetto  $p$  è definito directly density-reachable da un altro oggetto  $q$  se:

- 1)  $p \in \epsilon\text{-vicinato}(q)$
- 2)  $q$  è un CORE OBJECT

### DENSITY REACHABLE

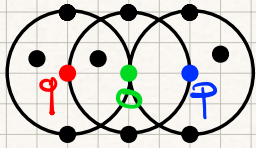
Un oggetto  $p$  è definito density-reachable da un altro oggetto  $q$  se esiste una catena di oggetti  $p_1 \dots p_n$  con  $p_1 = q$  e  $p_n = p$  tale per cui  $p_{i+1}$  è directly density-reachable da  $p_i$ .

Esempio:  $\epsilon = 1$  e  $\text{MinPoints} = 3$

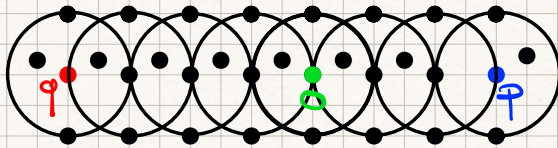


## DENSITY-CONNECTED

Un oggetto  $p$  è density-connected a un altro oggetto  $q$  se esiste un oggetto  $o$  tale per cui sia  $q$  che  $p$  sono density-reachable da  $o$ .



$\epsilon = 1$  e  $\text{MinPoints} = 3$





# DBSCAN

L'acronimo sta per **Density-Based Spatial Clustering of Application with Noise**.

L'idea alla base è semplice, un cluster è definito come un massimo insieme di **density-connected objects**.

Ogni qual volta trova un oggetto che ha nel  $\epsilon$ -vicinato più di minpoints un nuovo cluster viene generato con quel core-object. Poi iterativamente DBSCAN colleziona gli oggetti directly density-reachable dal core-object trovato. Questo può far venire la necessità di fare il merge di qualche density-reachable cluster.

**Stop-Condition:** Non ci sono nuovi punti da aggiungere a un cluster.

## DEFINIZIONE DI CLUSTER

Un sottoinsieme  $C$  in  $D$  è un cluster se:

- 1)  $\forall o_1, o_2 \in C : o_1$  e  $o_2$  sono density-connected
- 2) Non esiste nessun oggetto  $o \in C$  e  $o' \in (D-C)$  tale per cui  $o$  e  $o'$  sono density-connected

**COMPLESSITÀ:**  $O(n^2)$

## PARAMETRI:

- 1)  $\epsilon$
- 2) Min Points
- 3)  $D$

## ALGORITHM

DBSCAN( $\epsilon$ , minPoints)

```
{
  while (Non ho visitato tutti gli oggetti)
  {
    p = oggetto Random non visitato.
    p = visitato
    if ( $\epsilon$ -vicinato(p).size()  $\geq$  minPoints)
    {
      C = nuovo Cluster
      C.aggiungi(p)
      N =  $\epsilon$ -vicinato(p)

      for ( $\forall p' \in N$ )
      {
        if (p' non è stato visitato)
        {
          p' = visitato
          if ( $\epsilon$ -vicinato(p').size()  $\geq$  minPoints)
            N.add( $\epsilon$ -vicinato(p'))
          if (p' non appartiene a un cluster)
            C.aggiungi(p')
        }
      }
      return C
    }
    else { p' = rumore }
  }
}
```



## CONTRO

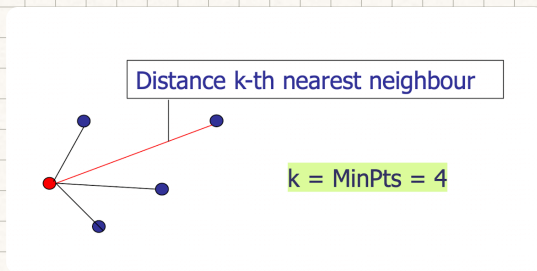
- 1) Molto sensibile ai parametri  $\epsilon$  e MinPoints. Se il dataset presenta densità diverse nello spazio DBSCAN potrebbe perdersi qualche cluster.

## Come determinare $\epsilon$ e MinPoints?

Usando un approccio euristico:

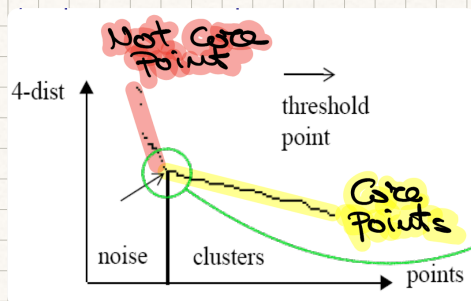
Possiamo settare MinPoints usando la distanza fra il punto preso in considerazione e il suo  $k$ -NN. Quindi diverso per ogni punto in  $D$ .

MinPoints sarà uguale al numero di punti all'interno del raggio dato da  $\epsilon = \text{dist}(p, k\text{-NN})$ , ovvero  $k$ .



## Possibile metodo euristico

- 1) Ordiniamo i punti in ordine decrescente di  $\text{dist}(p, k\text{-NN})$
- 2)  $\forall p \in D \quad \epsilon = \text{dist}(p, k\text{-NN}_p), \text{MinPoints} = k$ .  
Così tutti i punti con  $k$ -dist minore saranno CORE-OBJECT



- 3) Settiamo i nostri parametri con quelli inerenti al primo punto in cui la  $k$ -dist non cambia più drasticamente.

## PROBLEMA

Usare parametri generali per definire la densità, può non essere ottimale perché non tutto il DB ha la stessa densità.

### SOLUZIONE 1: HIERARCHICAL CLUSTER CON SINGLE-LINK

Questa soluzione ha però il difetto che due cluster separati dai pochi punti più vicini possono non essere propriamente separati. Il risultato è un dendrogram difficile da capire.

### SOLUZIONE 2: DENSITY-BASED PARTITION ALGORITHM

Il problema è la difficoltà di settare i parametri adeguatamente.

### SOLUZIONE BUONA:

Usare un algoritmo che produce un ordine speciale dei punti del dataset, basato sui livelli di densità presenti fra i vari potenziali cluster. Che è quello che fa OPTICS