

## PROXIMITY BASED

L'intuizione è che un punto molto lontano dagli altri è un outlier. Questo si basa sull'assunzione che gli outlier siano significativamente distanti dai punti normali.

## DISTANCE-BASED

Qui la prossimità è misurata con il vicinato, che se non è abbastanza popolato allora il punto in questione è un outlier.

Per ogni punto  $o$  si va a vedere il  $\epsilon$ -vicinato( $o$ ) dove  $\epsilon$  è una threshold in termini di raggio entro cui cercare vicini.

Un oggetto  $o$  è un outlier se:

$$\frac{|\{o' \mid \text{dist}(o, o') \leq \epsilon\}|}{\|D\|} \leq \pi$$

Ovvero se il numero di punti del suo  $\epsilon$ -vicinato in rapporto con la cardinalità di  $\|D\|$  è sotto una threshold  $\pi$  impostata dall'utente

## KNN

In alternativa possiamo calcolare la distanza fra  $o$  e il suo  $k$ -esimo vicino più vicino, dove  $k = \lceil \pi \cdot \|D\| \rceil$ . Se  $\text{dist}(o, o_k) > \tau$  allora outlier.

## Nested Loop algorithm

```
for ( $\forall o_i \in D$ )  
{  
  Calcola_distanza( $o_i$ , tutti gli altri punti)  
  Conta gli oggetti nel  $\epsilon$ -vicinato( $o_i$ )  
  if ( $|\epsilon\text{-vicinato}(o_i)| \geq \pi \cdot M$ )  
    continue  
  else  
     $o_i = \text{outlier}$   
}
```

## Grid-based algorithm

L'approccio di prima è dispendioso perché calcola tutte le distanze fra gli oggetti e il resto del DB. Perché non confrontare gruppi di oggetti? Ecco cosa fa l'algoritmo Grid-based.

Consiste nel dividere lo spazio dei dati in una griglia dove ogni cella ha diagonale pari a  $\epsilon/2$ .

### Tipi di celle:

1) LIVELLO 1: Presa una cella  $C$  di riferimento, sono tutte le celle distanti una cella da  $C$ . Formalmente  $\forall x \in C$  e  $\forall y \in C_{liv1} : \text{dist}(x, y) \leq \epsilon$

2) LIVELLO 2: Presa una cella  $C$  di riferimento, sono tutte le celle distanti più di 2 celle da  $C$ . Formalmente  $\forall x \in C$  e  $\forall y \in C_{liv2} : \text{dist}(x, y) > \epsilon$



## Regole di pruning

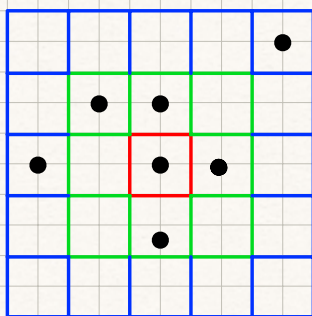
Sia " $a$ " il numero di punti in  $C$ , " $b_1$ " in  $C_{liv1}$  e " $b_2$ " in  $C_{liv2}$ .

1) REGOLA LIVELLO 1: Se  $a + b_1 \geq \pi n$  allora tutti gli oggetti  $\circ \in C$  NON SONO OUTLIER con  $\tau$  e  $\pi$  specificati. Questo perché tutti gli oggetti in  $C_{liv1}$  sono dentro l' $\tau$ -vicinato( $\circ$ ), e ci sono almeno  $\lceil \pi \tau \rceil$  vicini

2) REGOLA LIVELLO 2: Se  $a + b_1 + b_2 < \pi n + 1$  allora tutti gli oggetti  $\circ \in C$  SONO OUTLIER con  $\tau$  e  $\pi$  parametri dati, perché  $\forall \circ \quad |\tau\text{-vicinato}(\circ)| < \lceil \pi \tau \rceil$ .

Per la cella che soddisfa una delle due regole possiamo copiarla se nella cella sono presenti outlier.

## Esempio



Prendo a riferimento  $C$ ,  $C_{liv1}$  e  $C_{liv2}$

con  $\pi = 1/7$  e  $n = 7$

Per il Liv 1:

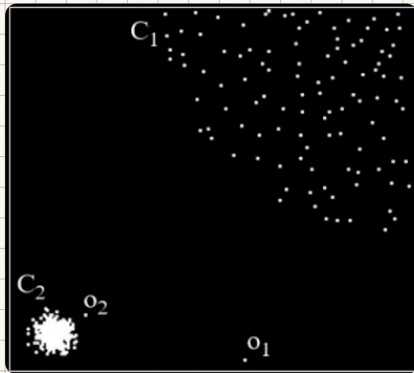
$$1 + 4 \geq 1 \quad \checkmark$$

Per il Liv 2

$$1 + 4 + 2 < 1 + 1 \quad \times$$

## LOCAL DISTANCE-BASED OUTLIER DETECT.

Ora abbiamo vere e proprii con le diverse densità che cluster diversi possono presentare nello stesso dataset, e che possono portare a non riconoscere qualche outlier



Esempio

Il punto  $o_2$  ha densità simile al cluster  $C_1$  se paragonato a  $C_2$ .

Definiamo:

1)  $N_k(x_i)$  i  $k$ -NN di  $x_i$

2)  $D_k(x_i)$  la distanza media fra  $x_i$  e  $k$ -NN

$$D_k(x_i) = \frac{1}{k} \sum_{x_j \in N_k(x_i)} |x_i - x_j|$$

Definiamo ora l'outlierness:

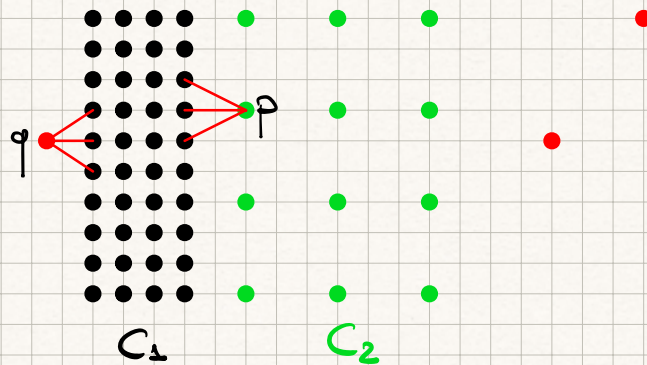
$$O_k(x_i) = \frac{D_k(x_i)}{\frac{1}{k} \sum_{j \in N_k(x_i)} D_k(x_j)}$$

Se  $O_k(x_i) > 1$  allora i vicini di  $x_i$  sono molto lontani



## Problema con l'outlieriness

Quando diversi cluster sono vicini fra loro,  $O_k(x_i)$  può dare risultati non attesi:



In questo caso  $O_3(p) > O_3(q)$  perché i punti verdi non fanno parte del 3-vicinato di  $q$  perché più distanti rispetto ai punti del cluster  $C_1$ .

# DENSITY-BASED OUTLIER DETECTION

L'intuizione è che la densità vicino a un punto outlier sarà molto diversa rispetto alle zone con punti normali.

Il metodo tende a paragonare la densità nel vicinato dell'oggetto sotto analisi con la densità degli oggetti vicini.

Definiamo:

1) **K-distance**:  $\text{dist}_k(o, K\text{-NN})$

2) **K-distance neighborhood of o**:  $N_k(o)$

$$N_k(o) = \{o' \mid o' \in D, \text{dist}(o, o') \leq \text{dist}_k(o, K\text{-NN})\}$$

$N_k(o)$  è l'insieme dei K-NN vicini di o.

Possono essere più di K considerando punti equidistanti da o

3) **REACHABILITY DISTANCE  $o'$  TO  $o$** :

$$\text{reachdist}(o \leftarrow o') = \max \{ \text{dist}_k(o), \text{dist}(o, o') \}$$

Quindi per i K-NN si prende in considerazione  $\text{dist}_k(o)$

**Proprietà:**  $\text{reachdist}(x \rightarrow y) \neq \text{reachdist}(y \rightarrow x)$

4) **LOCAL REACHABILITY DISTANCE OF  $o$** :

$$\text{Lrd}(o) = \frac{\|N_k(o)\|}{\sum_{o' \in N_k(o)} \text{reachdist}(o' \leftarrow o)}$$



# LOCAL OUTLIER FACTOR

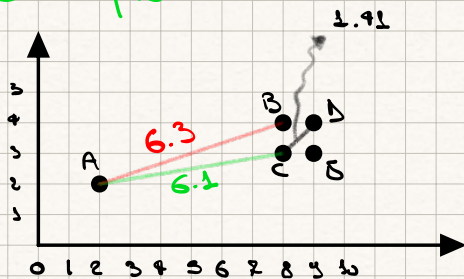
Per ogni oggetto possiamo calcolare questo fattore come:

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{reachdist(o, o')}{reachdist(o)} }{\|N_k(o)\|}$$

Questa formula dipende dal confronto fra  $reachdist(o)$  con  $reachdist(o')$  per tutto il vicinato, più il rapporto si alza, più LOF è alto.

Se  $reachdist(o) < reachdist(o') \Rightarrow LOF$  è alto

Esempio



Per  $K=2$ :

$$N_k(A) = \{B, C\}$$

$$dist_k(A) = dist(A, B) = 6.3$$

Calcoliamo ora:

$$reachdist(A \leftarrow B) = \max(6.3; 6.3) = 6.3$$

$$reachdist(A \leftarrow C) = \max(6.3; 6.1) = 6.3$$

Dato che servono, dobbiamo calcolare la reachability distance per il vicinato dei vicini di A.

$$N_k(B) = \{C, D, E\}$$

$$N_k(C) = \{B, D, E\}$$

$$dist_k(B) = dist(B, E) = \sqrt{2} = 1.41$$

$$dist_k(C) = dist(C, D) = \sqrt{2} = 1.41$$

Per B:

$$\text{reachdist}(B \leftarrow C) = 1.41$$

$$\text{reachdist}(B \leftarrow D) = 1.41$$

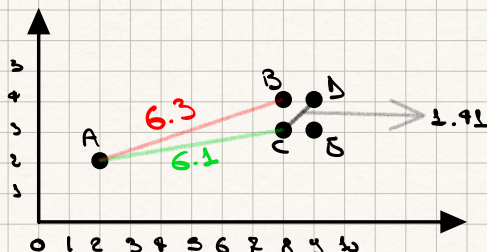
$$\text{reachdist}(B \leftarrow \delta) = 1.41$$

Per C:

$$\text{reachdist}(C \leftarrow B) = 1.41$$

$$\text{reachdist}(C \leftarrow D) = 1.41$$

$$\text{reachdist}(C \leftarrow \delta) = 1.41$$



Calcoliamo anche l'inverso per  $B \leftarrow A$  e  $C \leftarrow A$

$$\text{reachdist}(B \leftarrow A) = \max(1.41, 6.3) = 6.3$$

$$\text{reachdist}(C \leftarrow A) = \max(1.41, 6.1) = 6.1$$

Calcoliamo ora  $\text{ecc}_k(A)$ ,  $\text{ecc}_k(B)$  e  $\text{ecc}_k(C)$

$$\text{ecc}_k(A) = \frac{3}{6.3 + 6.1} = 0.24$$

$$\text{ecc}_k(B) = \frac{4}{1.41 + 1.41 + 1.41 + 1.41} = 0.71$$

$$\text{ecc}_k(C) = \frac{4}{1.41 + 1.41 + 1.41 + 1.41} = 0.71$$

Da cui ora possiamo calcolare il  $\text{WF}_k(A)$

$$\text{WF}_k(A) = \frac{\frac{0.71}{0.24} + \frac{0.71}{0.24}}{3} = 1.97$$



Problem 6F