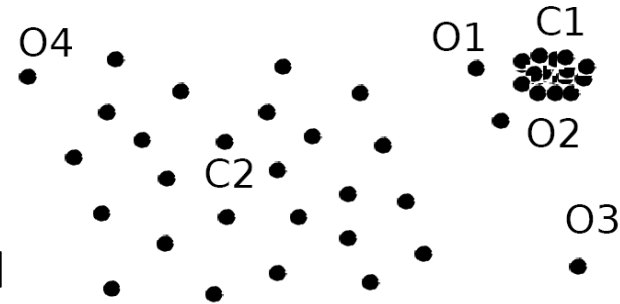


# Density-Based Outlier Detection

- Local outliers: Outliers comparing to their local neighborhoods, instead of the global data distribution
- In Fig.,  $o_1$  and  $o_2$  are local outliers to  $C_1$ ,  $o_3$  is a global outlier, but  $o_4$  is not an outlier. However, proximity-based clustering cannot find  $o_1$  and  $o_2$  are outlier (e.g., comparing with  $O_4$ ).



- Intuition (density-based outlier detection): The density around **an outlier** object is **significantly different from** the density around its neighbors
- Method: Use the relative density of an object against its neighbors as the indicator of the degree of the object being outliers
- *k-distance* of an object  $o$ ,  $\text{dist}_k(o)$ : distance between  $o$  and its  $k$ -th NN
- *k-distance neighborhood* of  $o$ ,  $N_k(o) = \{o' \mid o' \text{ in } D, \text{dist}(o, o') \leq \text{dist}_k(o)\}$ 
  - $N_k(o)$  could be bigger than  $k$  since multiple objects may have identical distance to  $o$

# Local Outlier Factor: LOF

- Reachability distance from  $o'$  to  $o$ :

$$reachdist_k(o \leftarrow o') = \max\{dist_k(o), dist(o, o')\}$$

– where  $k$  is a user-specified parameter

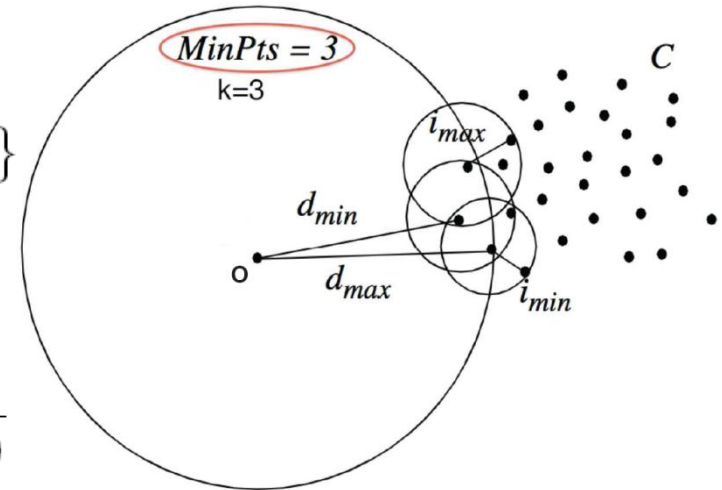
- Local reachability density of  $o$ :

$$lrd_k(o) = \frac{\|N_k(o)\|}{\sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)}$$

- LOF (Local outlier factor) of an object  $o$  is the average of the ratio of local reachability of  $o$  and those of  $o$ 's  $k$ -nearest neighbors

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{lrd_k(o')}{lrd_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)$$

- The lower the local reachability density of  $o$ , and the higher the local reachability density of the  $k$ NN of  $o$ , the higher LOF
- This captures a local outlier whose local density is relatively low comparing to the local densities of its  $k$ NN



# LOF(Local Outlier Factor) Example

- Consider the following 4 data points:

$a(0, 0)$ ,  $b(0, 1)$ ,  $c(1, 1)$ ,  $d(3, 0)$

Calculate the LOF for each point and show the top 1 outlier, set  $k = 2$  and use Manhattan Distance.

# Step 1: calculate all the distances between each two data points

- There are 4 data points:

$a(0, 0)$ ,  $b(0, 1)$ ,  $c(1, 1)$ ,  $d(3, 0)$

(Manhattan Distance here)

$$\text{dist}(a, b) = 1$$

$$\text{dist}(a, c) = 2$$

$$\text{dist}(a, d) = 3$$

$$\text{dist}(b, c) = 1$$

$$\text{dist}(b, d) = 3+1=4$$

$$\text{dist}(c, d) = 2+1=3$$

## Step 2: calculate all the $\text{dist}_2(o)$

- **$\text{dist}_k(o)$ : distance between  $o$  and its  $k$ -th NN(  $k$ -th nearest neighbor)**

**$\text{dist}_2(a) = \text{dist}(a, c) = 2$**  ( *$c$  is the 2<sup>nd</sup> nearest neighbor*)

**$\text{dist}_2(b) = \text{dist}(b, a) = 1$**  ( *$a/c$  is the 2<sup>nd</sup> nearest neighbor*)

**$\text{dist}_2(c) = \text{dist}(c, a) = 2$**  ( *$a$  is the 2<sup>nd</sup> nearest neighbor*)

**$\text{dist}_2(d) = \text{dist}(d, a) = 3$**  ( *$a/c$  is the 2<sup>nd</sup> nearest neighbor*)

## Step 3: calculate all the $N_k(o)$

- *k-distance neighborhood* of  $o$ ,  $N_k(o) = \{o' \mid o' \text{ in } D, \text{dist}(o, o') \leq \text{dist}_k(o)\}$

$$N_2(a) = \{b, c\}$$

$$N_2(b) = \{a, c\}$$

$$N_2(c) = \{b, a\}$$

$$N_2(d) = \{a, c\}$$

## Step 4: calculate all the $lrd_k(o)$

- $lrd_k(o)$ : Local Reachability Density of  $o$

$$lrd_k(o) = \frac{\|N_k(o)\|}{\sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)}$$

$$reachdist_k(o \leftarrow o') = \max\{dist_k(o), dist(o, o')\}$$

$\|N_k(o)\|$  means the number of objects in  $N_k(o)$ ,

For example:  $\|N_2(a)\| = \|\{b, c\}\| = 2$

$$lrd_k(a) = \frac{\|N_2(a)\|}{reachdist_2(b \leftarrow a) + reachdist_2(c \leftarrow a)}$$

## Step 4: calculate all the $\text{lrd}_k(o)$

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

$$\text{reachdist}_2(b \leftarrow a) = \max\{\text{dist}_2(b), \text{dist}(b, a)\}$$

$$= \max\{1, 1\} = 1$$

$$\text{reachdist}_2(c \leftarrow a) = \max\{\text{dist}_2(c), \text{dist}(c, a)\}$$

$$= \max\{2, 2\} = 2$$

Thus,  $\text{lrd}_2(a)$

$$= \frac{\|N_2(a)\|}{\text{reachdist}_2(b \leftarrow a) + \text{reachdist}_2(c \leftarrow a)} = 2/(1+2) = 0.667$$



## Step 4: calculate all the $\text{Ird}_k(o)$

Similarly,

$$\text{Ird}_2(b) = \frac{\|N_2(b)\|}{\text{reachdist}_2(a \leftarrow b) + \text{reachdist}_2(c \leftarrow b)} = 2/(2+2) = 0.5$$

$$\text{Ird}_2(c) = \frac{\|N_2(c)\|}{\text{reachdist}_2(b \leftarrow c) + \text{reachdist}_2(a \leftarrow c)} = 2/(1+2) = 0.667$$

$$\text{Ird}_2(d) = \frac{\|N_2(b)\|}{\text{reachdist}_2(a \leftarrow d) + \text{reachdist}_2(c \leftarrow d)} = 2/(3+3) = 0.33$$

# Step 5: calculate all the $LOF_k(o)$

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{lrd_k(o')}{lrd_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)$$

$$LOF_2(a) =$$

$$\begin{aligned} & (lrd_2(b) + lrd_2(c)) * (reachdist_2(b \leftarrow a) + reachdist_2(c \leftarrow a)) \\ & = (0.5 + 0.667) * (1 + 2) = 3.501 \end{aligned}$$

$$LOF_2(b) =$$

$$\begin{aligned} & (lrd_2(a) + lrd_2(c)) * (reachdist_2(a \leftarrow b) + reachdist_2(c \leftarrow b)) \\ & = (0.667 + 0.667) * (2 + 2) = 5.336 \end{aligned}$$

# Step 5: calculate all the $LOF_k(o)$

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{lrd_k(o')}{lrd_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)$$

$$LOF_2(c) =$$

$$\begin{aligned} & (lrd_2(b) + lrd_2(a)) * (reachdist_2(b \leftarrow c) + reachdist_2(a \leftarrow c)) \\ & = (0.5 + 0.667) * (1 + 2) = 3.501 \end{aligned}$$

$$LOF_2(d) =$$

$$\begin{aligned} & (lrd_2(a) + lrd_2(c)) * (reachdist_2(a \leftarrow d) + reachdist_2(c \leftarrow d)) \\ & = (0.667 + 0.667) * (3 + 3) = 8.004 \end{aligned}$$

## Step 6: Sort all the $\text{LOF}_k(o)$

The sorted order is:

$$\text{LOF}_2(\mathbf{d}) = 8.004$$

$$\text{LOF}_2(\mathbf{b}) = 5.336$$

$$\text{LOF}_2(\mathbf{a}) = 3.501$$

$$\text{LOF}_2(\mathbf{c}) = 3.501$$

**Obviously, top 1 outlier is point d.**