



程序代写代做 CS编程辅导



# Clustering and Density-based Anomaly Detection

WeChat: cstutorcs

Assignment Project Exam Help

COMP90073  
Email: tutorcs@163.com  
Security Analytics

QQ: 749389476  
Sarah Erfani, CIS

<https://tutorcs.com>  
Semester 2, 2021

# Outline

程序代写代做 CS编程辅导

- Anomaly detection with clustering
- Density-Based Spatial Clustering of Applications with Noise (DBSCAN)
- Local Outlier Factor (LOF)



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- **Advantages:**

- They can detect anomalies without requiring any *labelled* data.
- They work for many data types.
- Clusters can be regarded as summaries of the data.
- Once the clusters are obtained, clustering-based methods need only compare any object against the clusters to determine whether the object is an anomaly.
- Test process is typically fast and efficient because the number of clusters is usually small compared to the total number of objects small.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

- **Weakness:**

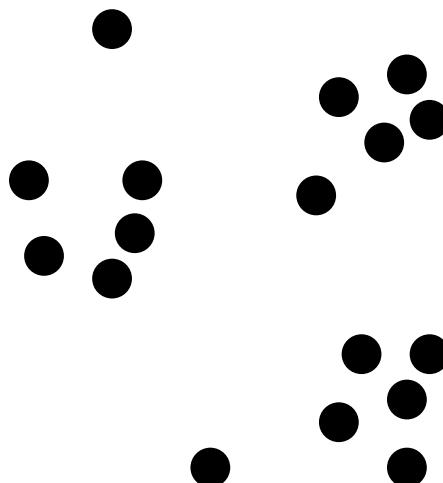
QQ: 749389476

- Their effectiveness depends highly on the clustering method used. Such methods may not be optimized for outlier detection.
- They are often costly for large data sets, which can serve as a bottleneck.

<https://tutorcs.com>

# Clustering-based Anomaly Detection

- $k$ -means clustering



程序代写代做 CS 编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Clustering-based Anomaly Detection

- $k$ -means clustering



WeChat: cstutorcs  
 $\xrightarrow{k=3}$

Assignment Project Exam Help

Email: tutorcs@163.com

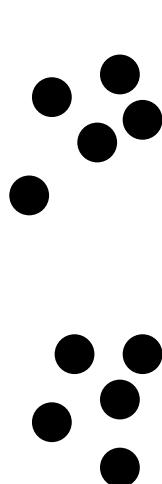
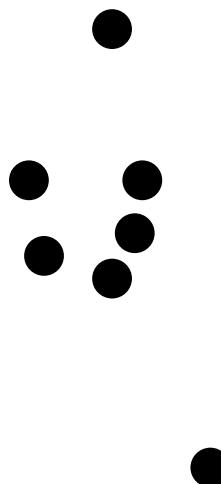
QQ: 749389476

<https://tutorcs.com>

+ Cluster center

# Clustering-based Anomaly Detection

- $k$ -means clustering



WeChat: cstutorcs  
 $\xrightarrow{k=3}$

Assignment Project Exam Help

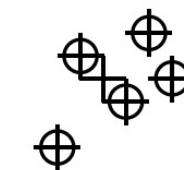
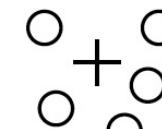
Email: tutorcs@163.com

QQ: 749389476

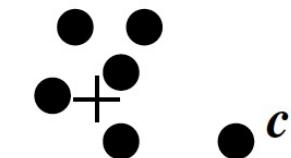
<https://tutorcs.com>

程序代写代做 CS 编程辅导

$a$



$b$



$c$

+ Cluster center

# Clustering-based Anomaly Detection

程序代写代做 CS编程辅导

- Assign an anomaly score to each object according to the distance between the object and the centre of closest cluster.



- Anomaly score( $p_j$ ) =  $\frac{dist(p_j, c_0)}{\frac{1}{n} \sum_i dist(p_i, c_0)}$

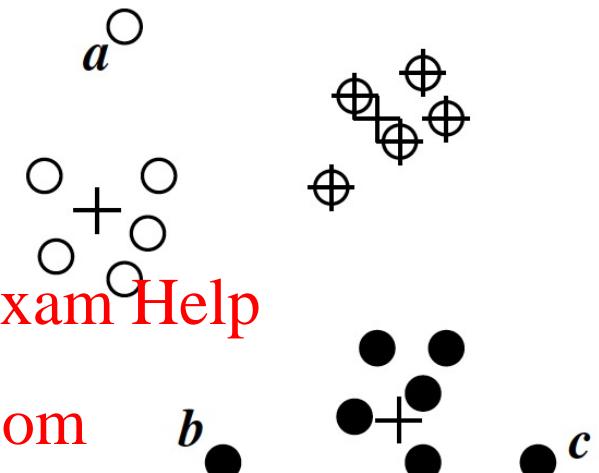
WeChat: cstutorcs

Assignment Project Exam Help

- Anomalies (**a**,**b**,**c**) are far from the clusters to which they are closest (with respect to the cluster centres).

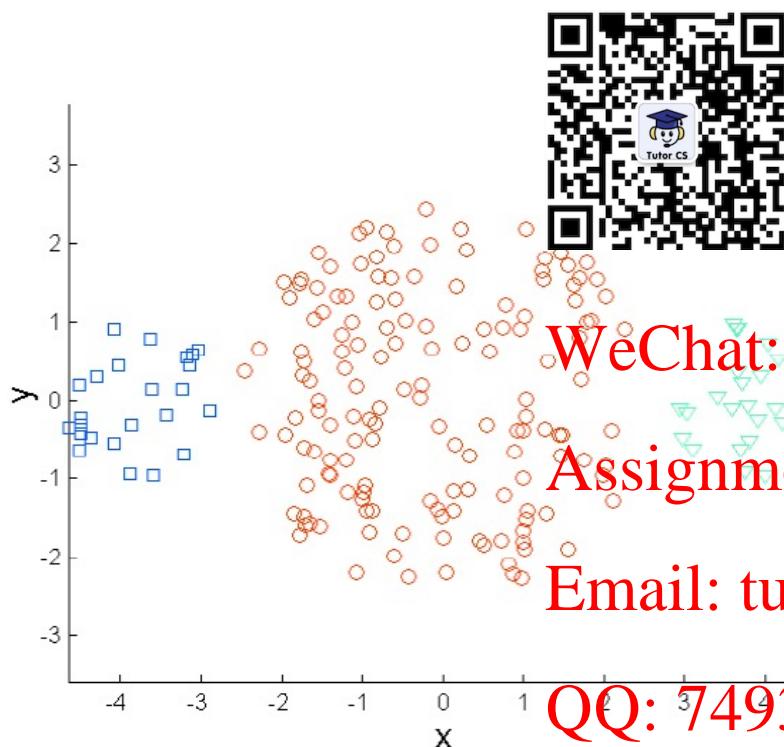
Email: tutorcs@163.com  
QQ: 749389476

<https://tutorcs.com>



# Limitations of *k*-means: Differing Size

程序代写代做 CS编程辅导



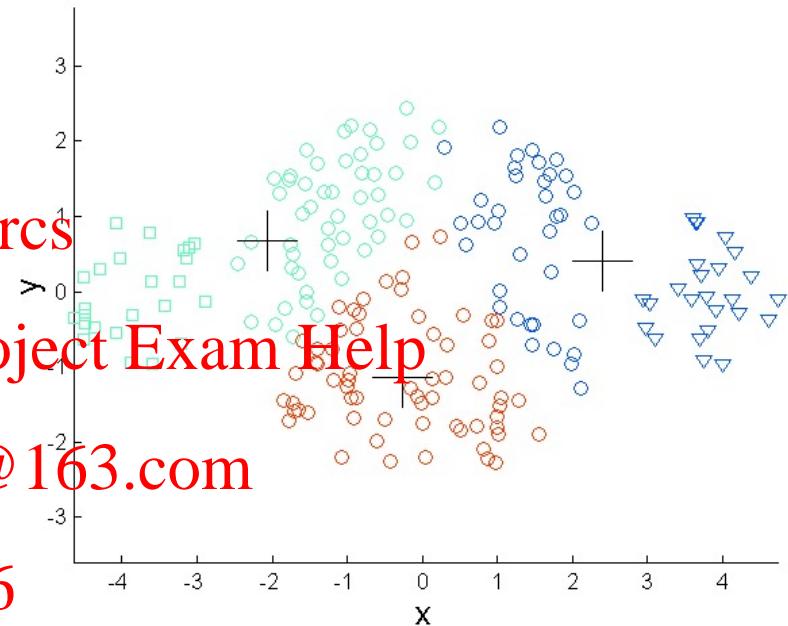
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

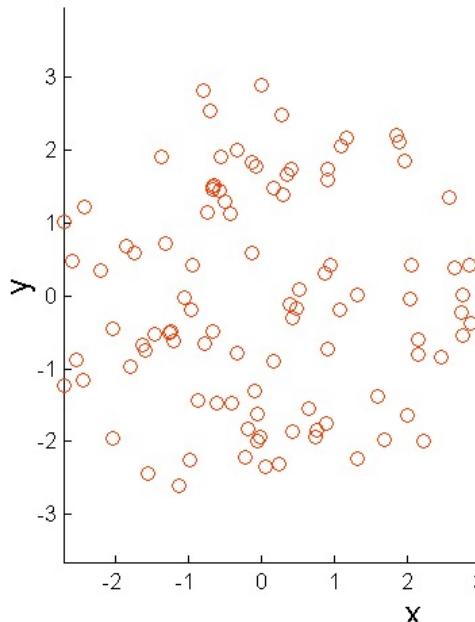
Original Points <https://tutorcs.com>



K-means (3 Clusters)

# Limitations of *k*-means: Differing Density

程序代写代做 CS编程辅导

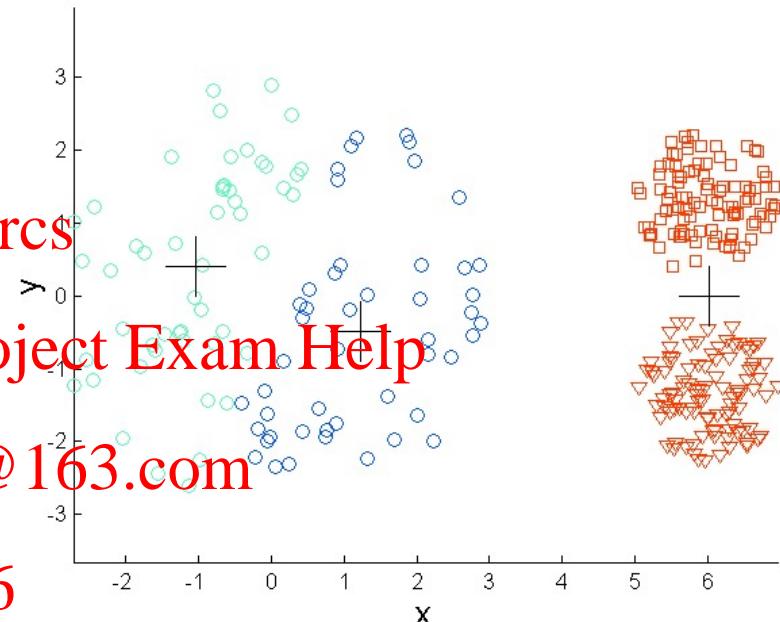


WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

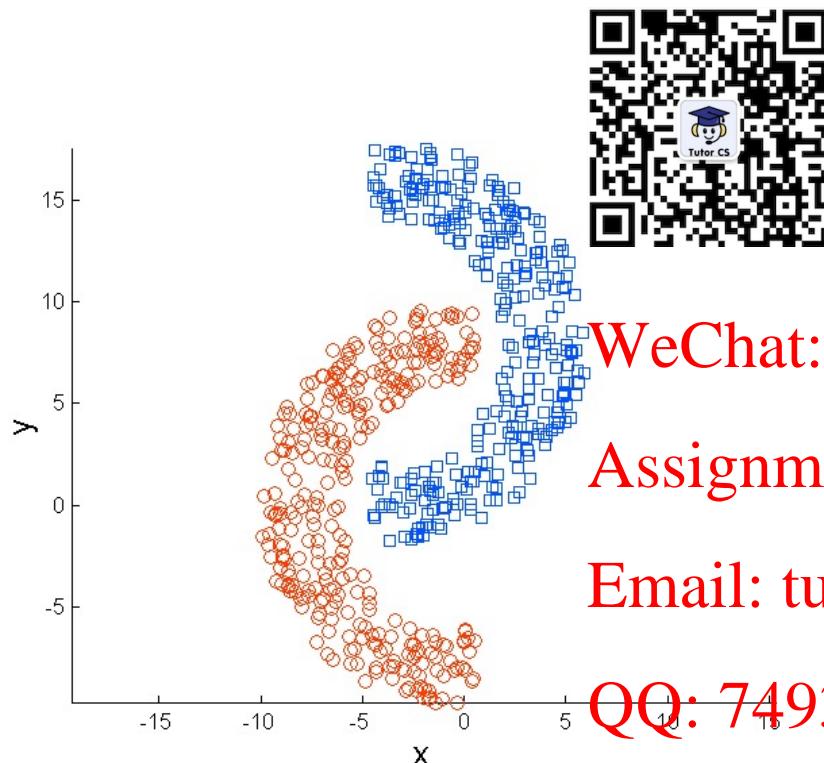


Original Points <https://tutorcs.com>

K-means (3 Clusters)

# Limitations of $k$ -means: Non Globular Shape

程序代写代做 CS编程辅导



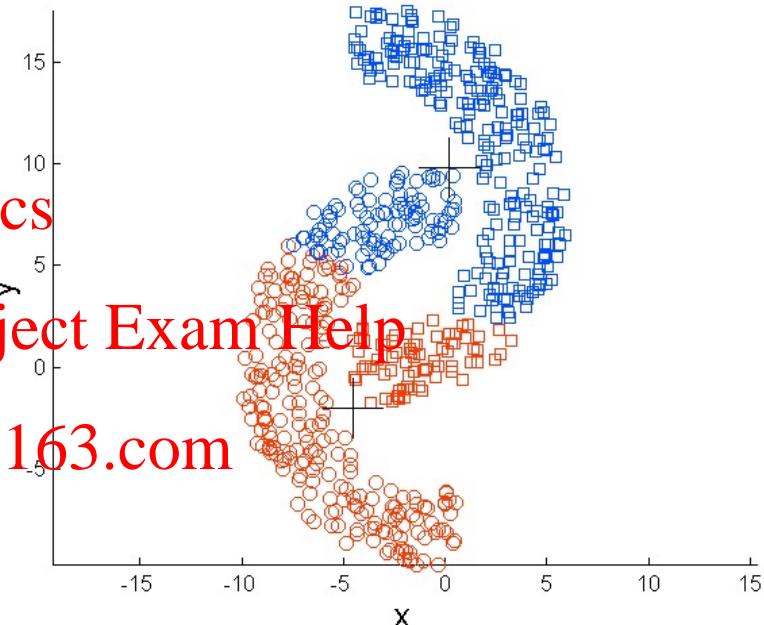
Original Points <https://tutorcs.com>

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476



K-means (2 Clusters)

# Density based Clustering

程序代写代做 CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

- **Density-based clustering:** Model clusters as dense regions in the data space, separated by sparse regions, which can discover clusters of non-spherical shape and avoid outliers.

<https://tutorcs.com>

程序代写代做 CS 编程辅导

- **Objective:** Identify dense regions, which can be measured by the number of objects close to a given point.  
— Finds core objects, then finds objects that have dense neighbourhoods. It connects core objects with their neighbourhoods to form dense regions as clusters.

WeChat: cstutorcs

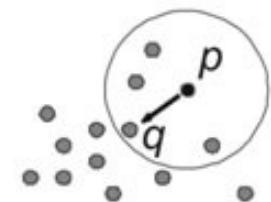
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Density at point p: Number of points within a circle of radius Eps

Dense Region: A cluster with radius Eps that contains at least MinPts points



Eps = 1cm  
MinPts = 4

Parameters:

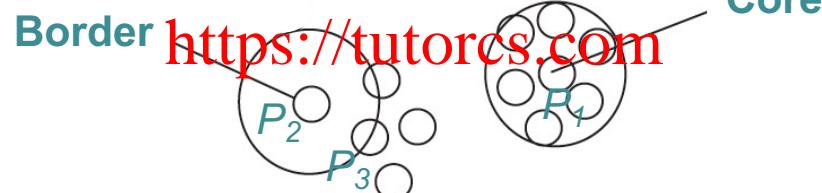
程序代写代做 CS编程辅导

DBSCAN defines different classes of points:

- **Core point:** A point with at least MinPts points within its Eps-neighbourhood (including *itself*). 
- **Border point:** A point with fewer points than MinPts in the Eps-neighbourhood, but is in the neighbourhood of a core point.
- **Anomaly (outlier) point:** A point which is neither core nor border.
- E.g., MinPts = 4

Email: [tutorcs@163.com](mailto:tutorcs@163.com)

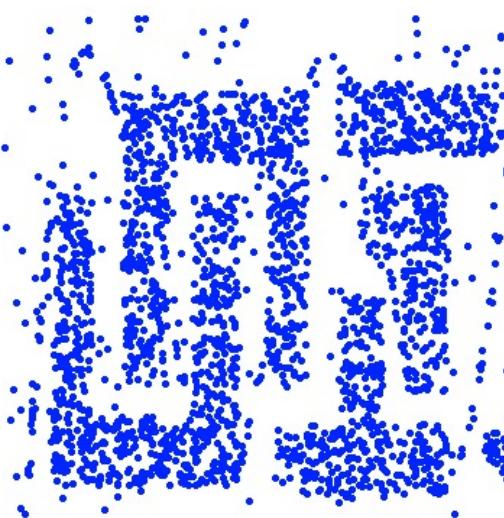
QQ: 749389476



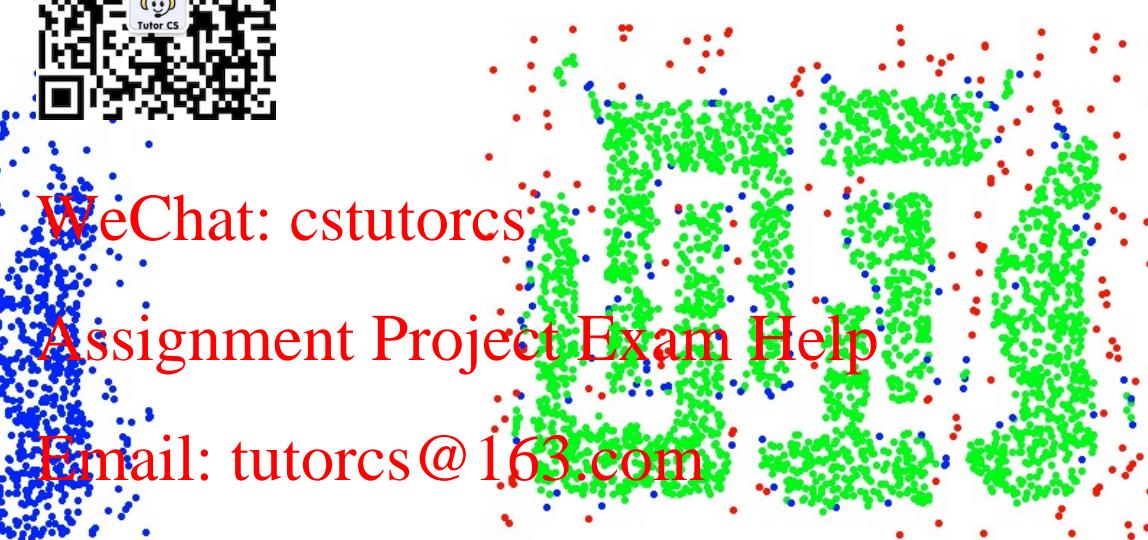
Border <https://tutorcs.com>

Core

程序代写代做 CS编程辅导



Original data



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

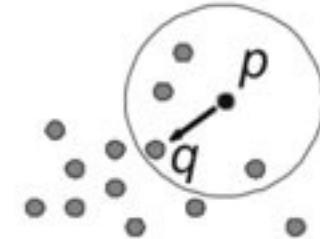
QQ: 749389476

<https://tutorcs.com>

Point types: core, border  
and anomaly

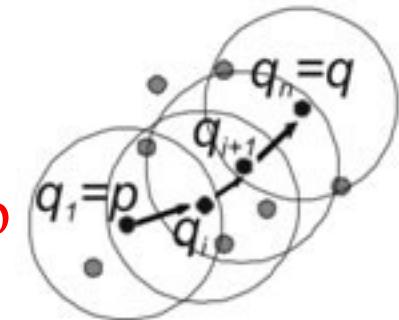
程序代写代做 CS 编程辅导

- **Directly Density-reachable:** Point  $q$  is directly density-reachable from  $p$  (w.r.t. Eps and MinPts) if  $p$  is a *core point*, and  $q$  is within the  $\text{Eps}$ -neighbourhood of  $p$ .



- **(Indirectly) Density-reachable:** Point  $q$  is density-reachable from  $p$  (w.r.t. Eps and MinPts) if there is a chain of points  $q_1, \dots, q_n$ ,  $q_1 = p$ ,  $q_n = q$ , such that  $q_{i+1}$  is directly density-reachable from  $q_i$ .

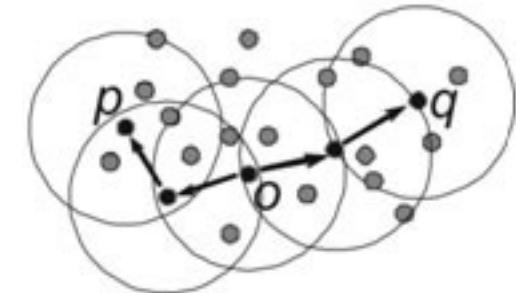
WeChat: tutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com



- **Density-connected:** Point  $q$  is density-connected to a point  $p$  (w.r.t. Eps and MinPts) if there is a point  $o$  such that both  $p$  and  $q$  are density reachable from  $o$  (w.r.t. Eps and MinPts).

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$



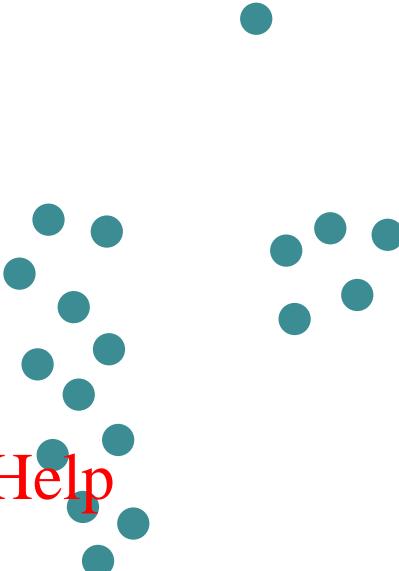
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$



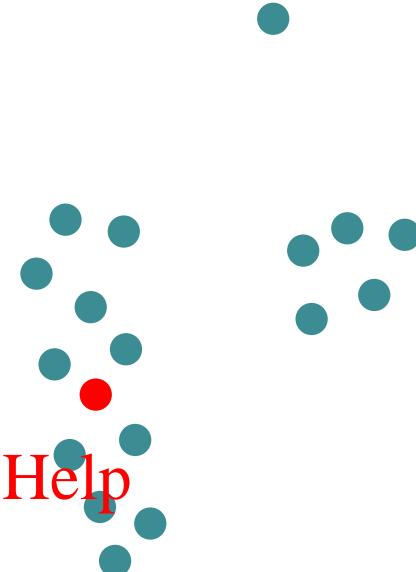
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

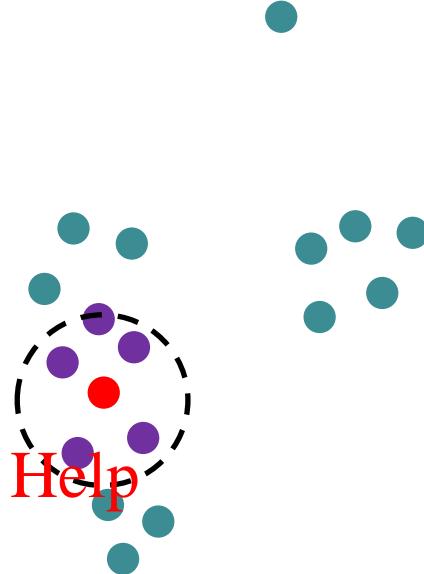
QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
- Retrieve all points density  $\text{DB}(p)$  from  $p$  w.r.t. Eps and MinPts (e.g.



WeChat: cstutorcs

Assignment Project Exam Help

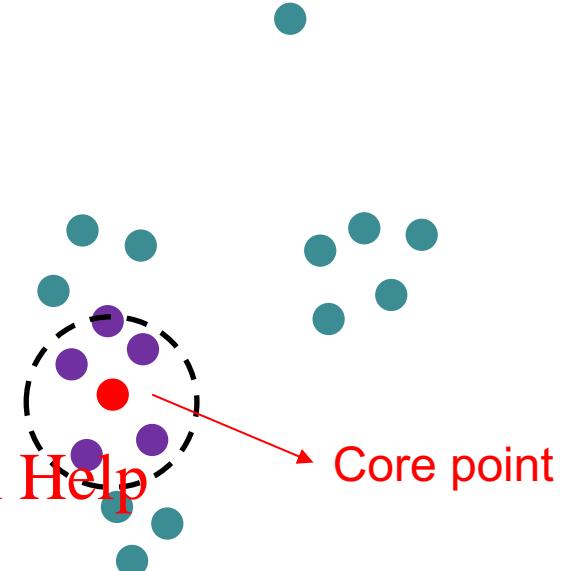
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
- Retrieve all points density  $\text{neighbor}(p)$  from  $p$  w.r.t. Eps and MinPts (e.g.
- If  $p$  is a core point, create



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
- Retrieve all points density  $\text{neighbor}(p)$  from  $p$  w.r.t. Eps and MinPts (e.g.
- If  $p$  is a core point, create



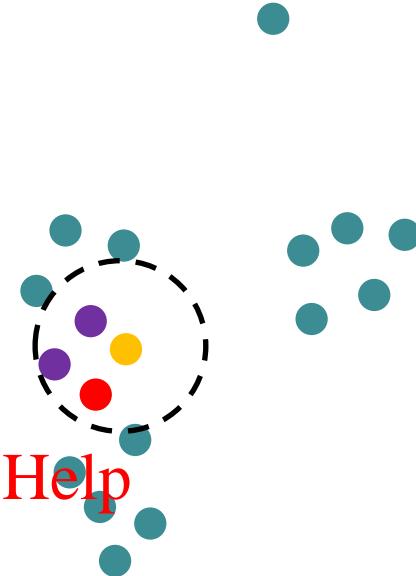
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



# DBSCAN Algorithm

程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
- Retrieve all points density-reachable from  $p$  w.r.t. Eps and MinPts (e.g.
- If  $p$  is a core point, create
- If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database



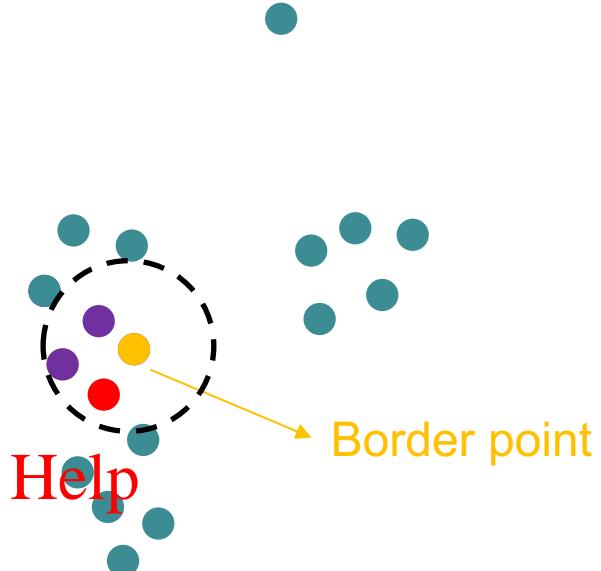
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
- Retrieve all points density-reachable from  $p$  w.r.t. Eps and MinPts
- If  $p$  is a core point, create a cluster
- If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database
- Repeat the above steps until all data points have been visited

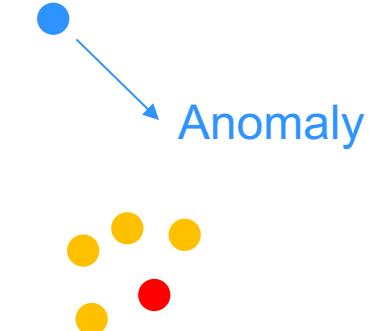
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



# DBSCAN Algorithm

程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
- Retrieve all points density-reachable from  $p$  w.r.t. Eps and MinPts
- If  $p$  is a core point, create a cluster
- If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database
- Repeat the above steps until all data points have been visited

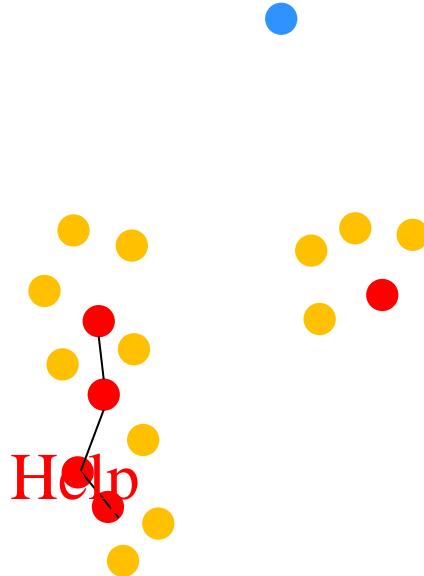
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

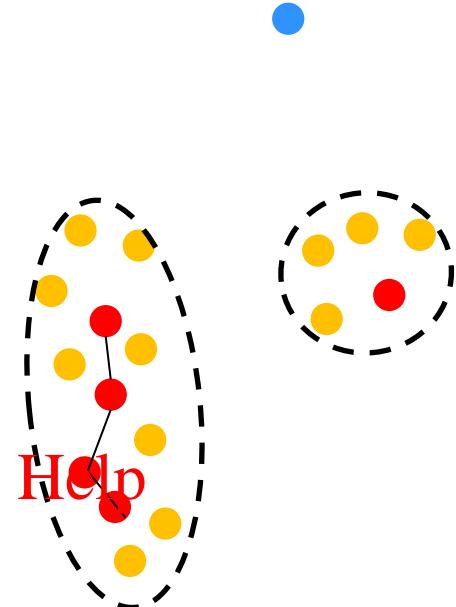
- Randomly select an unvisited object  $p$
- Retrieve all points density-reachable from  $p$  w.r.t. Eps and MinPts
- If  $p$  is a core point, create a cluster
- If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database
- Repeat the above steps until all data points have been visited

WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



# 程序代写代做 CS编程辅导

- Randomly select an unvisited object  $p$
  - Retrieve all points density-reachable from  $p$  w.r.t. Eps and MinT
  - If  $p$  is a core point, create
  - If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database
  - Repeat the above steps until all data points have been visited

WeChat: cstutorcs

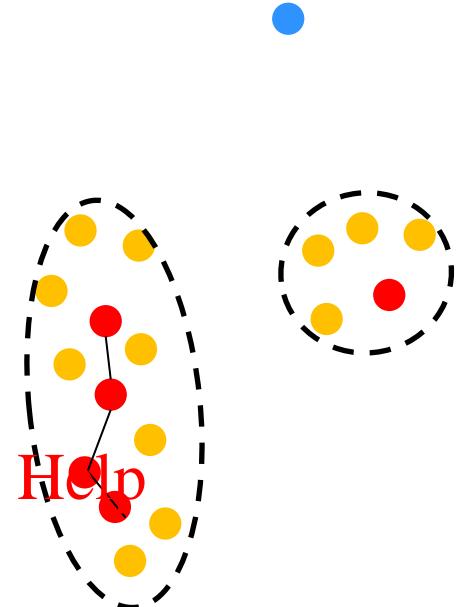
Assignment Project



# SCAN visits the next WeChat: cstutorcs

Assignment Project Exam Help until all data points

Email: tutorcs@163.com



**Computational Complexity**: QQ: 749389476

- $O(n^2)$ , where  $n$  is the number of samples.
  - If a spatial index is used,  $O(n \log n)$ .

# Determining Eps and MinPts

程序代写代做 CS编程辅导

- Idea is that for points in a cluster, their  $k^{\text{th}}$  nearest neighbour is roughly the same distance



- Noise points have the  $k^{\text{th}}$  neighbour at farther distance

WeChat: cstutorcs

- So, plot sorted distance of every point to its  $k^{\text{th}}$  nearest neighbour

Assignment Project Exam Help

- Find the distance  $d$  where there is a “knee” in the curve

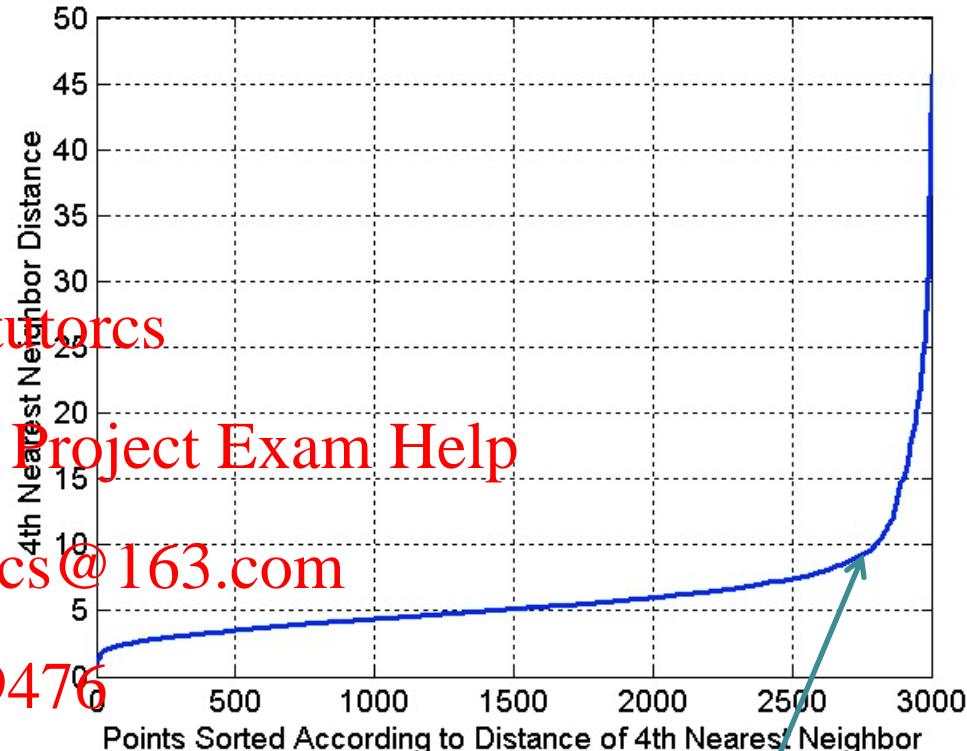
—  $\text{Eps} = d$ ,  $\text{MinPts} = k$

QQ: 749389476

- Demo:

<https://tutorcs.com>

<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>



- **Advantages:**

- Resistant to Noise
- Can handle clusters of different shapes and sizes

- **Disadvantages:**

- Varying densities



shapes and sizes



- Sensitive to parameter setting

QQ: 749389476



Eps = 5

Eps = 3.5

Eps = 3

- High-dimensional data

程序代写代做 CS编程辅导

- **Advantages:**

- They can detect anomalies without requiring any labelled data.
- They work for many different types of data.
- Clusters can be regarded as summaries of the data.
- Once the clusters are obtained, clustering-based methods need only compare any object against the clusters to determine whether the object is an anomaly.
- Test process is typically faster than other methods because the number of clusters is usually small compared to the total number of objects.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

- **Weakness:**

- Their effectiveness depends highly on the clustering method used. Such methods may not be optimized for anomaly detection.
- They are often costly for large data sets, which can serve as a bottleneck.



# Local Proximity-based Outliers

程序代写代做 CS编程辅导

In the following figure which of the following instances are anomalies?

- $o_1$ ?
- $o_2$ ?
- $o_3$ ?
- $o_4$ ?



程序代写代做 CS编程辅导

- **Objective:** Quantify the *relative density* about a particular data point.
- **Intuition:** The anomalies are more *isolated* compared to “normal” data points.
- LOF uses the relative density of an object against its neighbours to indicate the degree to which an object is an anomaly.



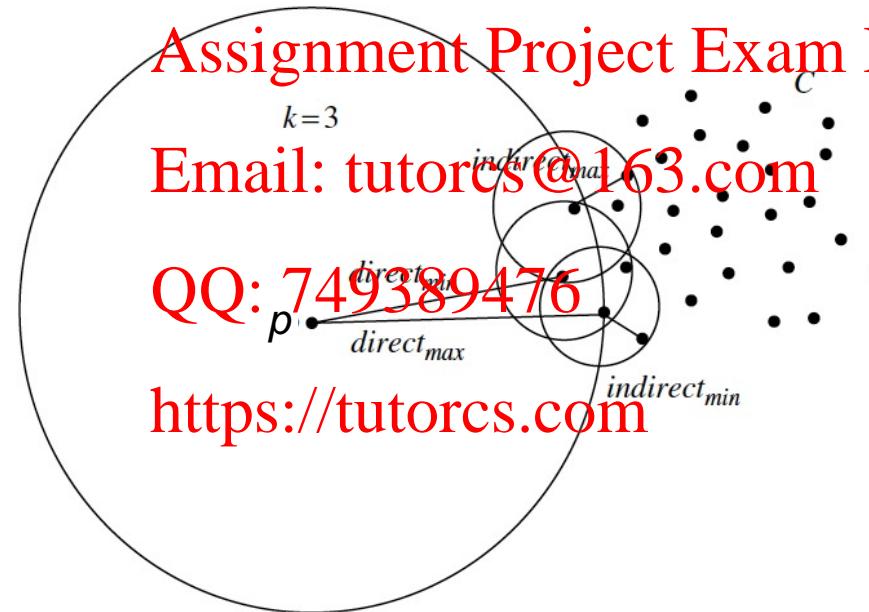
WeChat: cstutorcs

Assignment Project Exam Help

 $k=3$ 

Email: tutorcs@163.com

QQ: 749389476

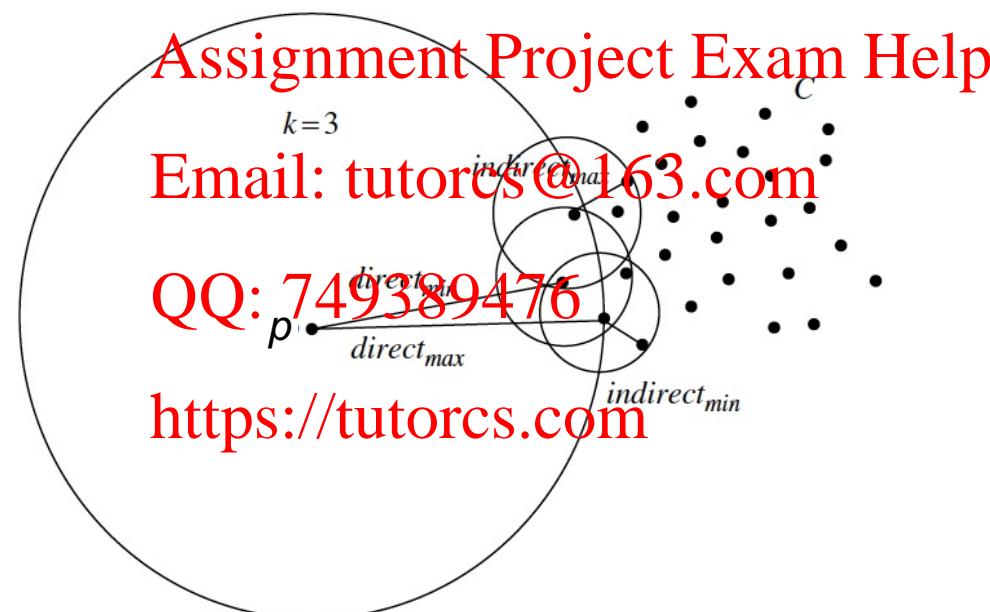
<https://tutorcs.com>

程序代写代做 CS编程辅导

- $kdist$ : distance between  $p$  and its  $k^{th}$  NN
- Meta-heuristic: The  $kdist$  notion of “volume”
- The more isolated a point is, the larger its  $kdist$



WeChat: cstutorcs



# Reachability Distance

程序代写代做 CS编程辅导

- Reachability Distance of  $p$  with respect to  $o$ :

$$reachdis(p, o) = \max\{kdist(o), dist(p, o)\}$$


Not symmetric

- Intuition: “Do your close neighbours see you as one of their close neighbours”

WeChat: cstutorcs

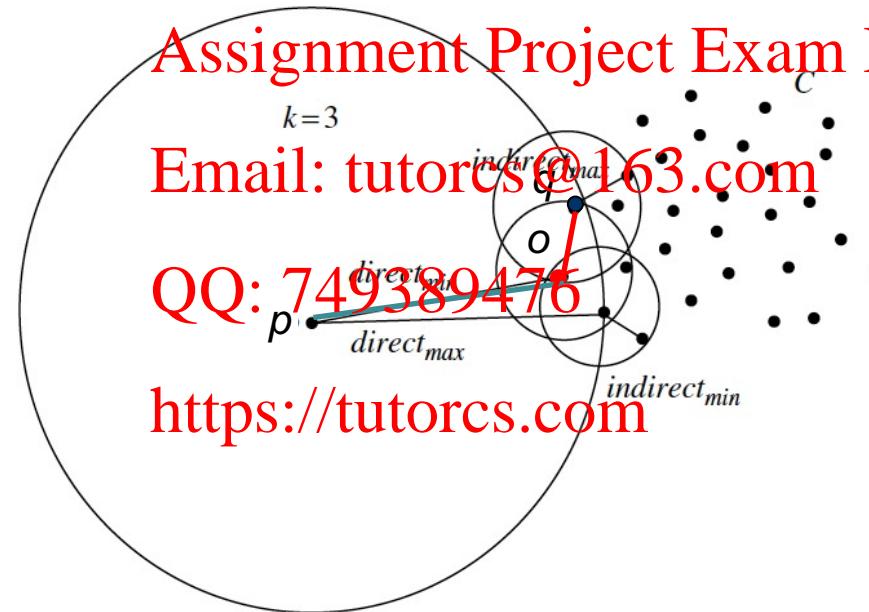
Assignment Project Exam Help

$k=3$

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

- Local Reachability Density of  $p$ :

$$lrd_k(p) = \left( \frac{1}{\sum_{o \in N_{(p,k)}} reachdist_k(p, o)} \right)^{-1}$$

WeChat: cstutorcs  
 k nearest neighbours of  $p$

Assignment Project Exam Help

- **Intuition:** How far we have to travel from our point to reach the next point or cluster of points.
  - The lower it is, the less dense it is, the longer we have to travel.

<https://tutorcs.com>

程序代写代做 CS编程辅导

- LOF of an object  $p$  is the average of the ratio of local reachability of  $p$  and those of  $o$ 's k-nearest neighbour
- The anomalies are coming from this dense area, so the ratio is higher for anomalies



$$LOF_k(p) = \frac{1}{k} \sum_{o \in N_{(p,k)}} \frac{lrd_k(o)}{lrd_k(p)}$$

WeChat: cstutorcs

Assignment Project Exam Help



Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- The lower the local reachability density of  $p$ , and the higher the local reachability density of the  $\text{neighbors}$  of  $p$ , the higher LOF



- $LOF_k(p) \sim 1$ : Comparable density to neighbours,
- $LOF_k(p) < 1$ : Higher density than neighbours
- $LOF_k(p) > 1$ : Lower density than neighbours

WeChat: cstutors

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Data set

LOFs

程序代写代做 CS编程辅导

- Consider the following 4 data points:

a(0, 0), b(0, 1), c(1,



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

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS 编程辅导

Step 1: Calculate all the distances between each two data points

- There are 4 data points:

a(0, 0), b(0, 1), c(1,



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

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

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

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

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

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

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS 编程辅导

Step 2: Calculate  $\text{dist}_k(o)$ , distance between o and its k-th NN (k-th nearest neighbour)

k=2:

$\text{dist}_2(a) = \text{dist}(a, c) = 2$  (c is the 2nd nearest neighbour)

$\text{dist}_2(b) = \text{dist}(b, a) = 1$  (a/c is the 2nd nearest neighbour)

$\text{dist}_2(c) = \text{dist}(c, a) = 2$  (a is the 2nd nearest neighbour)

$\text{dist}_2(d) = \text{dist}(d, a) = 3$  (a/c is the 2nd nearest neighbour)



Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Step 3: Calculate all the  $N_k(p)$ , k-distance neighborhood of  $p$ ,  $N_k(p) = \{p' | p' \text{ in } D, \text{dist}(p, p') \leq \text{dist}_k(p)\}$

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

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

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

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



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Step 4: Calculate all the  $lrd_k(p)$

程序代写代做 CS编程辅导

For example:



$$lrd_k(a) = \frac{\|N_2(a)\|}{reachdist(a, b) + reachdist(a, c)}$$

WeChat: cstutorcs

- $\|N_2(a)\| = \|\{b, c\}\| = 2$

Assignment Project Exam Help

- $reachdist(a, b) = \max\{dist_2(b)dist(b, a)\} = \max\{1, 1\} = 1$

Email: tutorcs@163.com

- $reachdist(a, c) = \max\{dist_2(c)dist(c, a)\} = \max\{2, 2\} = 2$

QQ: 749389476

$$lrd_k(a) = \frac{2}{1 + 2} = 0.67$$

程序代写代做 CS编程辅导

Step 4: Calculate all the  $lrd_k(p)$

Similarly,



- $lrd_k(b) = \frac{\|N_2(b)\|}{reachdist(b,a)+reachdist(b,c)} = \frac{2}{2+2} = 0.5$

Assignment Project Exam Help

- $lrd_k(c) = \frac{\|N_2(c)\|}{reachdist(c,b)+reachdist(c,a)} = \frac{1}{1+2} = 0.33$

Email: tutorcs@163.com  
QQ: 749389476

- $lrd_k(d) = \frac{\|N_2(d)\|}{reachdist(d,a)+reachdist(d,c)} = \frac{2}{3+3} = 0.33$

程序代写代做 CS编程辅导  
Step 5: calculate all the  $LOF_k(p)$



- $LOF_2(a) = (lrd_2(b) + lrd_2(c)) \times (reachdist_2(a, b) + reachdist_2(a, c)) = (0.5 + 0.67) \times (1 + 2) = 3.51$
- $LOF_2(b) = (lrd_2(a) + lrd_2(c)) \times (reachdist_2(b, a) + reachdist_2(b, c)) = (0.67 + 0.67) \times (2 + 2) = 5.36$

Assignment Project Exam Help

- $LOF_2(c) = (lrd_2(b) + lrd_2(a)) \times (reachdist_2(c, b) + reachdist_2(c, a)) = (0.5 + 0.67) \times (1 + 2) = 3.51$
- $LOF_2(d) = (lrd_2(a) + lrd_2(c)) \times (reachdist_2(d, a) + reachdist_2(d, c)) = (0.67 + 0.67) \times (3 + 3) = 8.04$

QQ: 749389476

<https://tutorcs.com>

Step 6: Sort all the  $LOF_k(p)$

程序代写代做 CS编程辅导

The sorted order is:

- $LOF_2(d) = 8.04$
- $LOF_2(d) = 5.36$
- $LOF_2(d) = 3.51$
- $LOF_2(d) = 3.51$



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Obviously, top 1 anomaly is point d

程序代写代做 CS编程辅导

- LOF captures a *local anomaly* whose local density is relatively low comparing to the local densities of its neighbors
- Outputs a *scoring* (assigns a  $\text{LOF}_{\text{point}}$  value to each point)
- Choice of  $k$  specifies the reference set
- Originally implements a local approach (resolution depends on the user's choice for  $k$ )

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

1) Find clusters in a data set (using k-means)

2) Sort them according to size.

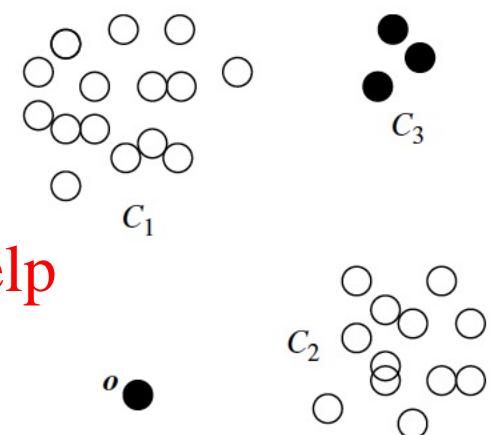
- Any cluster that contains at least a percentage (e.g., 90%) of the data set is considered a “large cluster.” The remaining clusters are referred to as “small clusters.”



WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com  
QQ: 749389476

3) To each data point, assign a cluster-based local

outlier factor (CBLOF), which is computed as the  
*product of the cluster's size and the similarity*  
*between the point and the cluster.*



- For a point belonging to a small cluster, its CBLOF is calculated as the product of the size of the small cluster and the similarity between the point and the closest large cluster.

程序代写代做 CS编程辅导

- What are the advantages of clustering for anomaly detection?
- How distance and density based clustering perform differently?
- How to identify local anomalies?
- How to identify group anomalies?

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

Next: Anomaly Detection in Evolving Data Streams

<https://tutorcs.com>

## References

程序代写代做 CS 编程辅导

1. Jiawei Han, Micheline Kamber, Jian Pei, "Data Mining: Concepts and Techniques", 3<sup>rd</sup> ed, 2011. Ch 10. Clusters 10.4 and 12



2. Martin Ester , Hans-Peter Kriegel , Jörg Sander , Xiaowei Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise.", KDD, 1996.

3. Density-Based Clustering

<http://www.cse.buffalo.edu/faculty/azhang/cse601/density-based.ppt>

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>