



程序代写代做 CS编程辅导



Andromeda Detection in Evolving Data Streams

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

- Introduction to data streams
- Windowing techniques
- HS-Trees
- Incremental LOF (iLOF)



WeChat: cstutorcs

Assignment Project Exam Help

— Memory-efficient iLOF (MiLOF)
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Data Streams

程序代写代做 CS编程辅导

Data stream is a sequence of data points, which is *continues*, *unbounded*, and *nonstationary*.



- **Streamlining Analysis**

- Large volume of data
- Short/real-time response
- Limited memory
- Energy/communication constraints

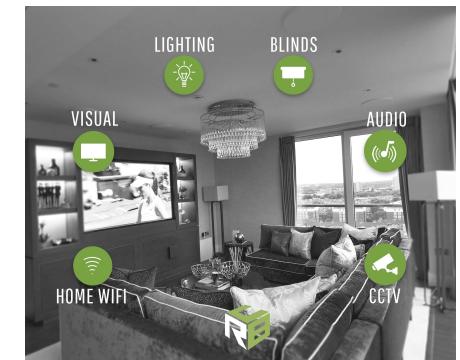
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Batch Learning vs. Incremental Learning

Batch Learning: Data points are stored until they can be analysed at the end of a monitoring period. Batch learning methods

- Can be computationally expensive
- Their accuracy is heavily dependent on a good choice for the training period and the quality of the training data
- Cannot be applied in *streaming environments*, where the measurements arrive as a continuous stream of data
- Cannot be used in *resource constraint devices* to buffer all the measurements
- Cannot identify anomalous points as *they occur*
- Cannot adapt to changes in the environment (e.g., drift)



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Incremental Learning: Data points are (usually) analysed once and there is *no need to buffer the data*. Incremental methods

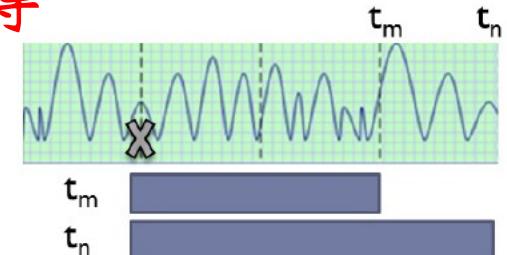
<https://tutorcs.com>

- Start with a set of initial parameters for the selected model and they becomes more accurate as more data arrives

Different Windowing Techniques for Data Streams

程序代写代做 CS编程辅导

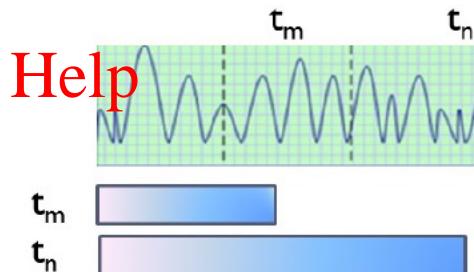
- **Landmark windows:** A fixed point in the data stream is defined as a landmark. Processing is done over data points between the landmark and the present data point.



- **Damped windows:** A weight is assigned to each data point in such a way that the old data points are given smaller weights. Therefore, the most recent data points are used with higher weights.

Assignment Project Exam Help

Email: tutorcs@163.com



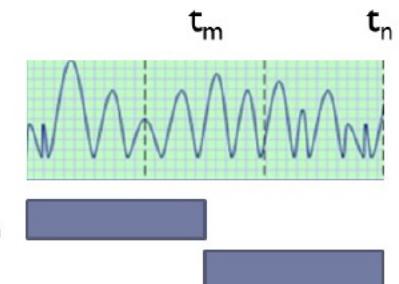
QQ: 749389476

<https://tutorcs.com>

Different Windowing Techniques for Data Streams

程序代写代做 CS编程辅导

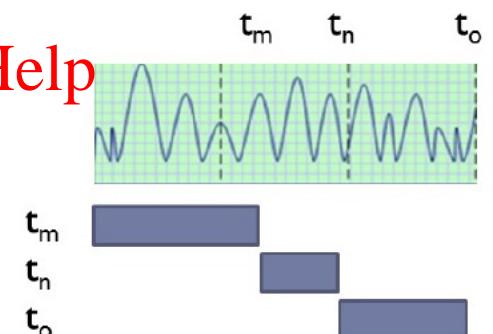
- **Sliding windows:** A sliding window size w is considered in this technique. It processes the last w data points in the stream. Older data points are discarded.



WeChat: cstutorcs

- **Adaptive windows:** The window size w would change as the data stream evolves. In this technique, the more the data points evolve, the smaller w becomes. In contrast, if data points remain constant, the value of w increases.

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 fast one-class anomaly detector for evolving data streams.

- A random tree model
- Builds tree structure without



- Detects anomalies in one pass

WeChat: cstutorcs
Assignment Project Exam Help

- Adapts to distribution changes by regular model updates

Email: tutorcs@163.com

QQ: 749389476

- Requires constant amount of memory $O(t^{2^h})$

<https://tutorcs.com>

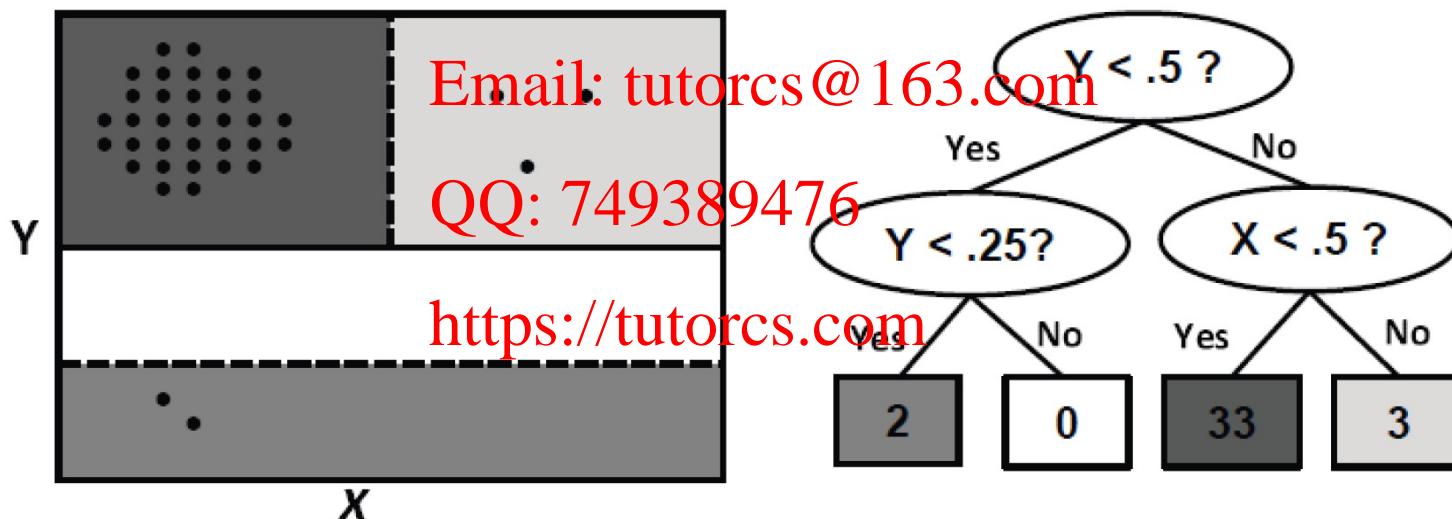
t : number of trees, h : depth of tree, ψ : window size

程序代写代做 CS 编程辅导

An HS-Tree is a full binary tree, which all leaves are at the same depth h .

- Randomly select an attribute d
- Bisects the space into two half-spaces, using the mid-point of d (assume that attributes' ranges are normalized to $[0, 1]$)
- Continue expansion until the maximum depth h of all nodes is reached.
- Employs mass as a measure to rank anomalies.

Assignment Project Exam Help



Separating the Regions: HS-Trees

程序代写代做 CS编程辅导



Ranking by Mass

程序代写代做 CS编程辅导



程序代写代做 CS编程辅导

- Divide data stream into fixed-size windows: W_1, W_2, \dots, W_n
- Each window is a fixed number of sequenced data instances
- **Initial Learning:** Train model M_1 using instances in W_1
- **Subsequent Learning and Anomaly Scoring**



WeChat: cstutorcs

For each window W_k (where $k > 1$)

Assignment Project Exam Help

Train model M_k using instances in W_k

Email: tutorcs@163.com

Test instances in W_k using model M_{k-1}

QQ: 749389476

Next window

<https://tutorcs.com>

- Let window size = 3
- Initial stage
- W1: reference window
 - Train HS-Trees and Classify



程序代写代做 CS编程辅导

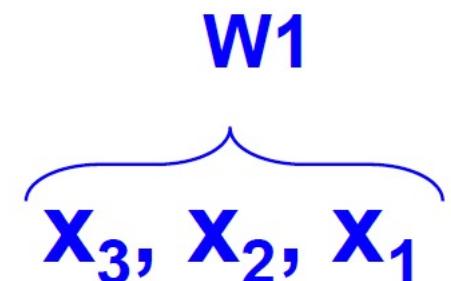
WeChat: cstutorcs

Assignment Project Exam Help

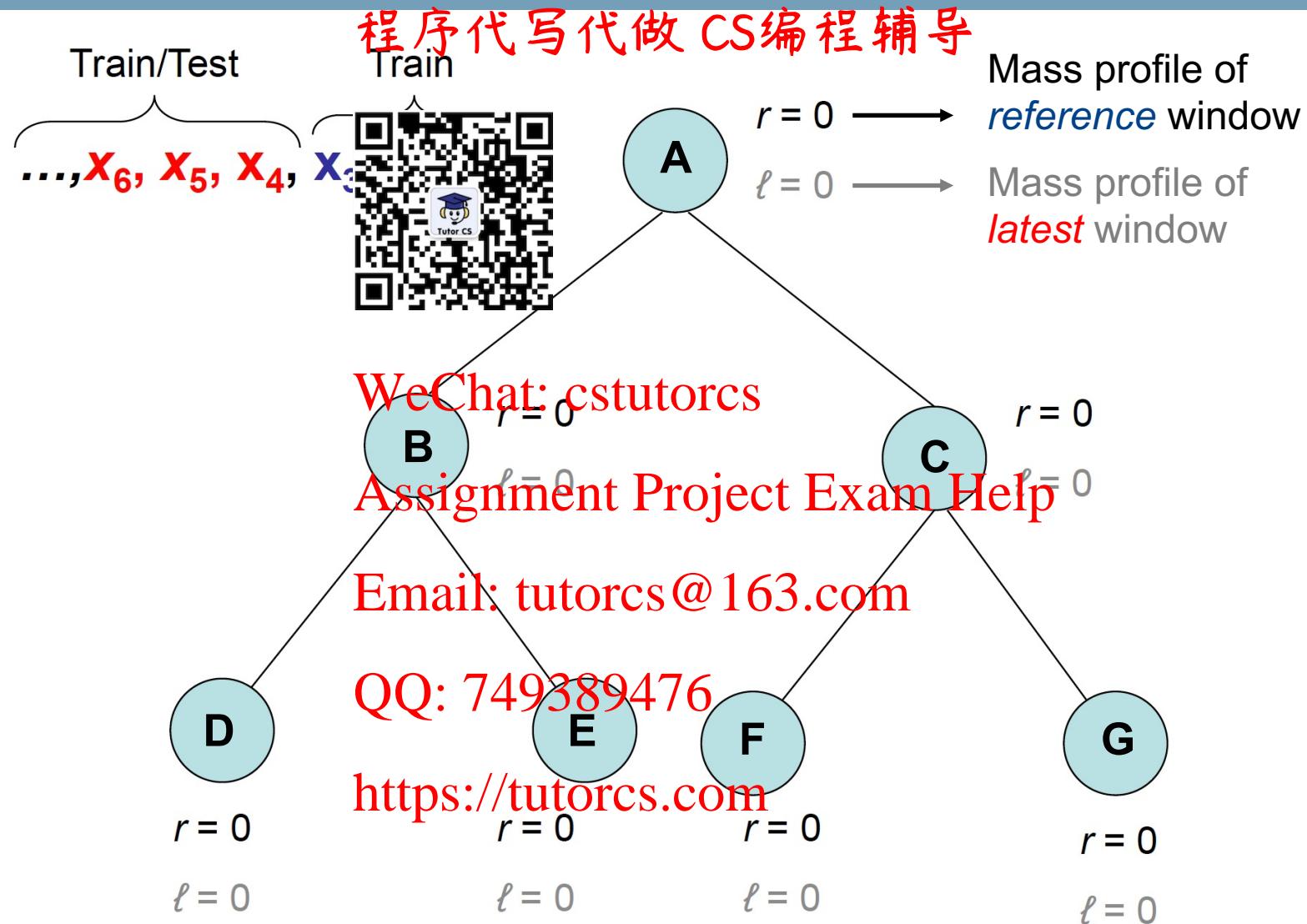
Email: tutorcs@163.com

QQ: 749389476

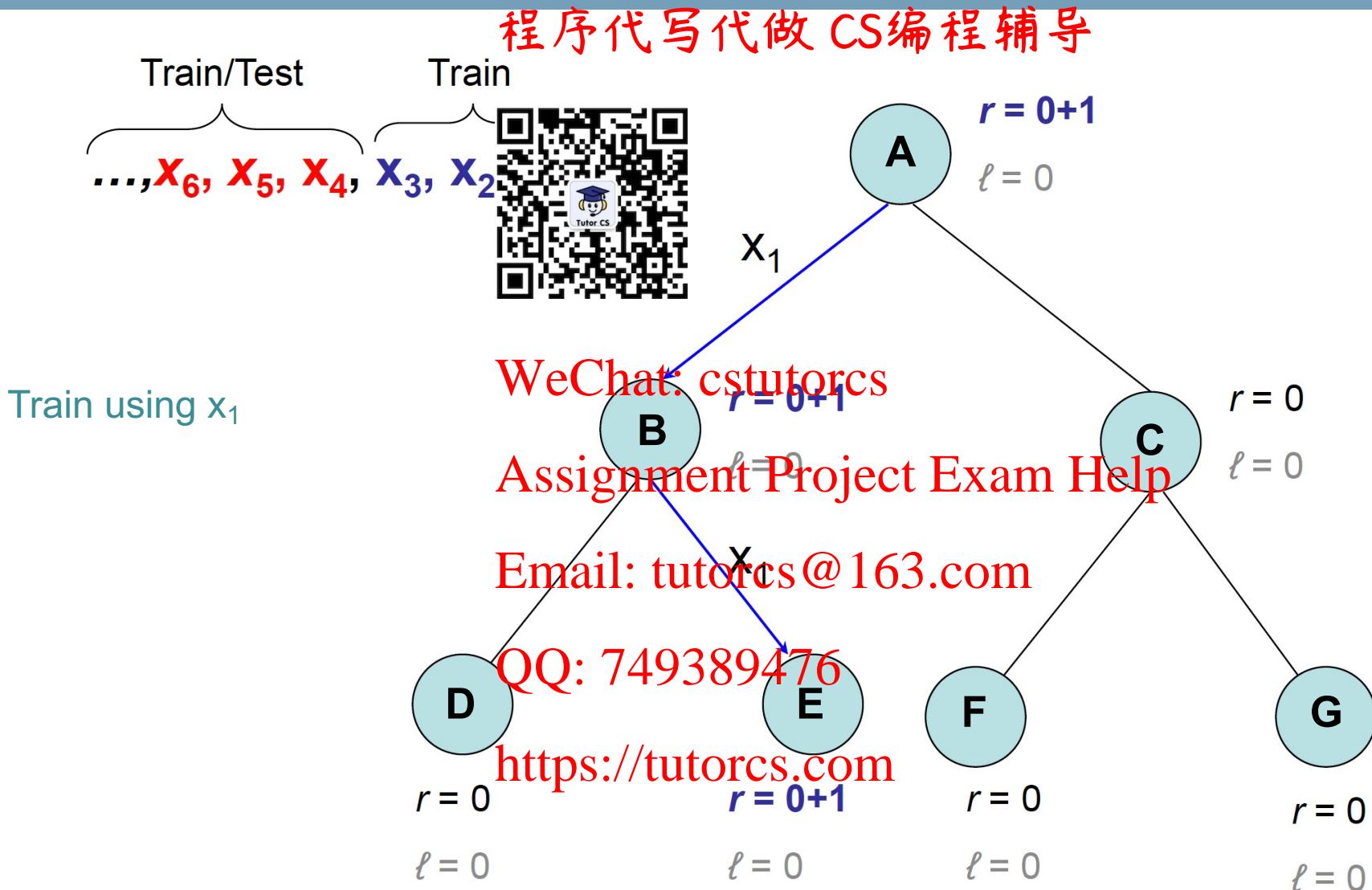
<https://tutorcs.com>



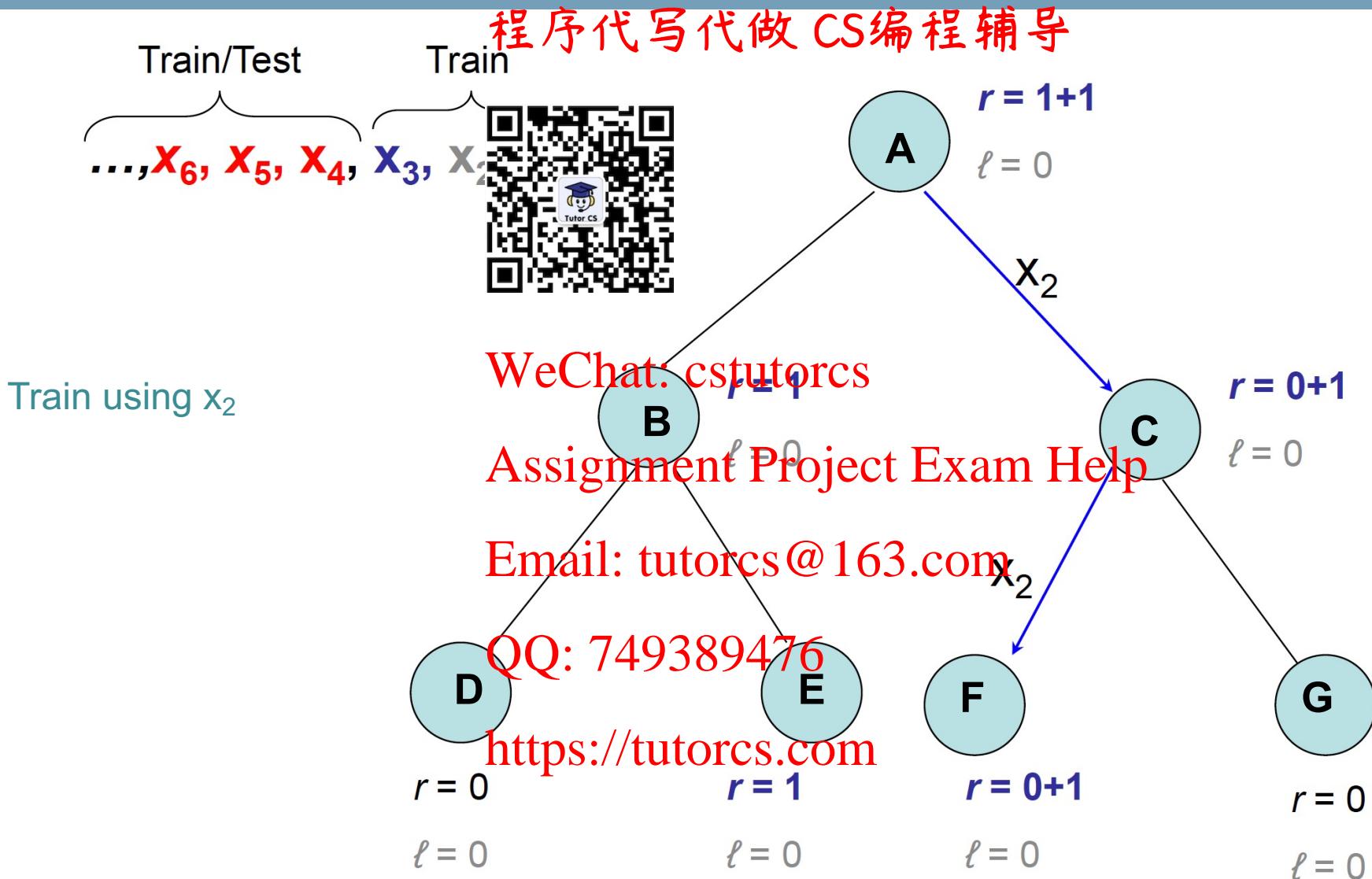
Streaming HS-Trees – Example



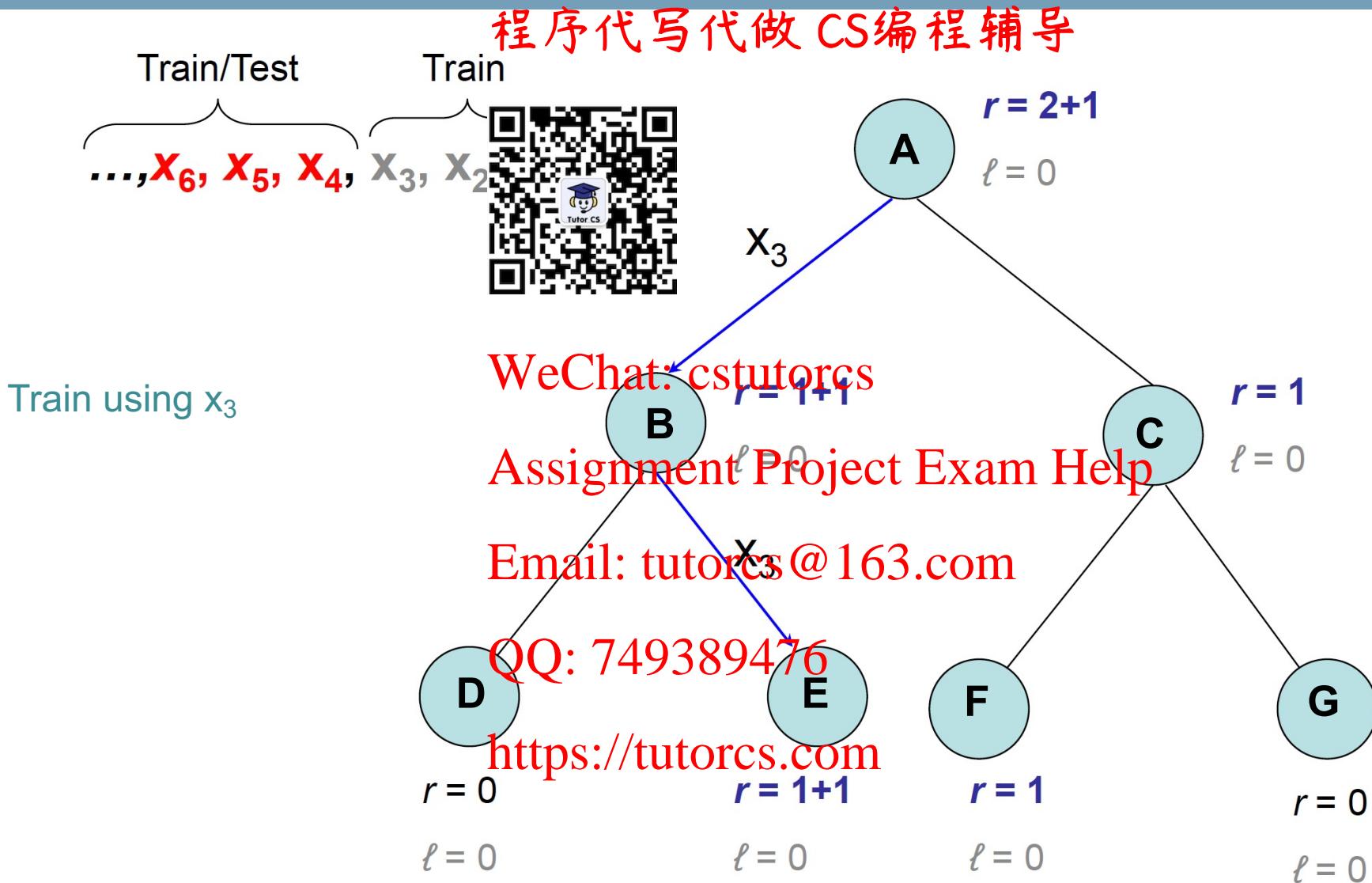
Streaming HS-Trees – Example



Streaming HS-Trees – Example

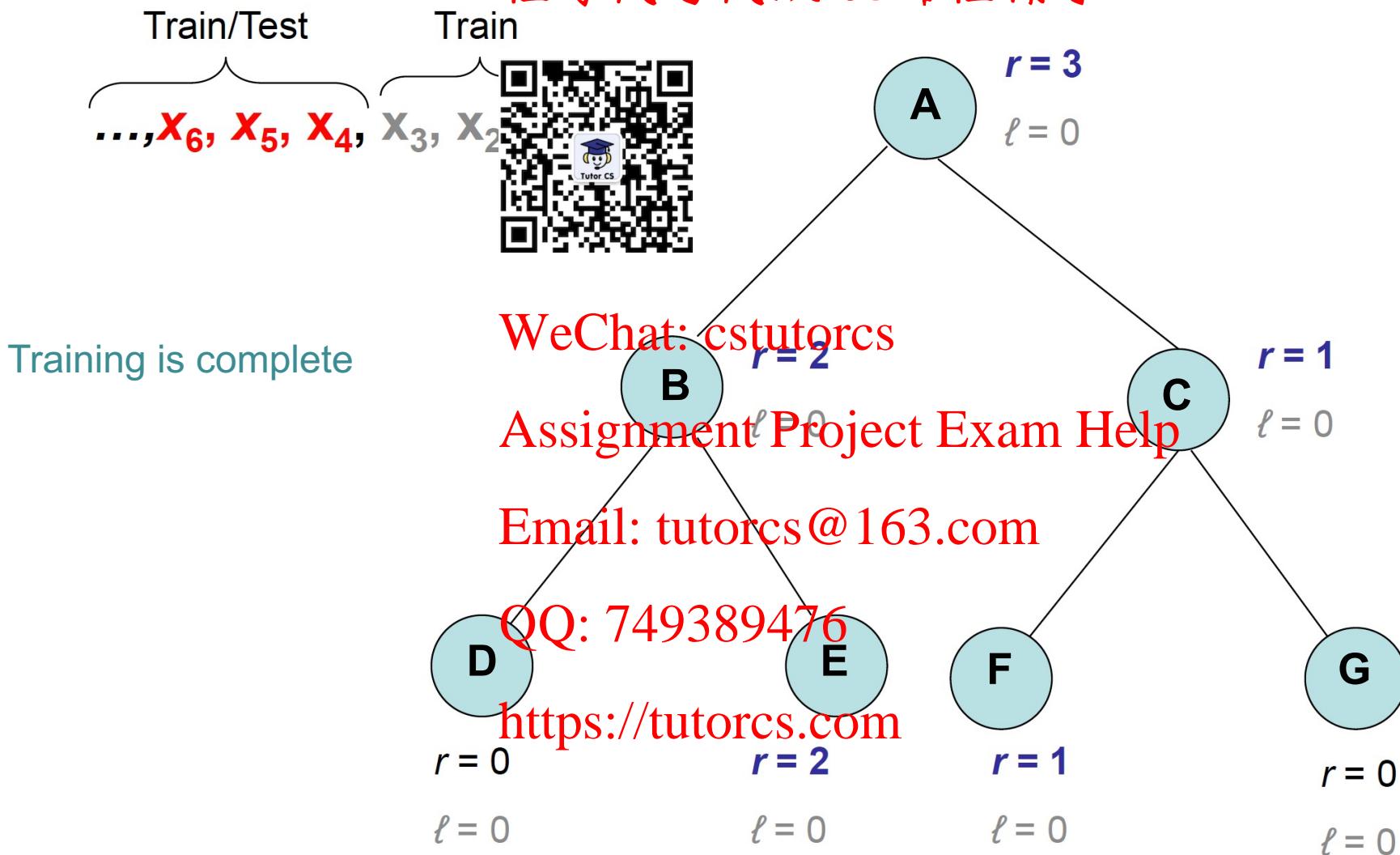


Streaming HS-Trees – Example



Streaming HS-Trees – Example

程序代写代做 CS 编程辅导



程序代写代做 CS 编程辅导

- Let window size = 3
- Initial stage
- W1: reference window
 - Train HS-Trees and mass r
- W2: latest window
 - Instances in W2 for training HS-Trees (mass ℓ)
 - Instances in W2 for testing HS-Trees (mass r)



WeChat: **TutorCS_tutors**
Assignment Project Exam Help

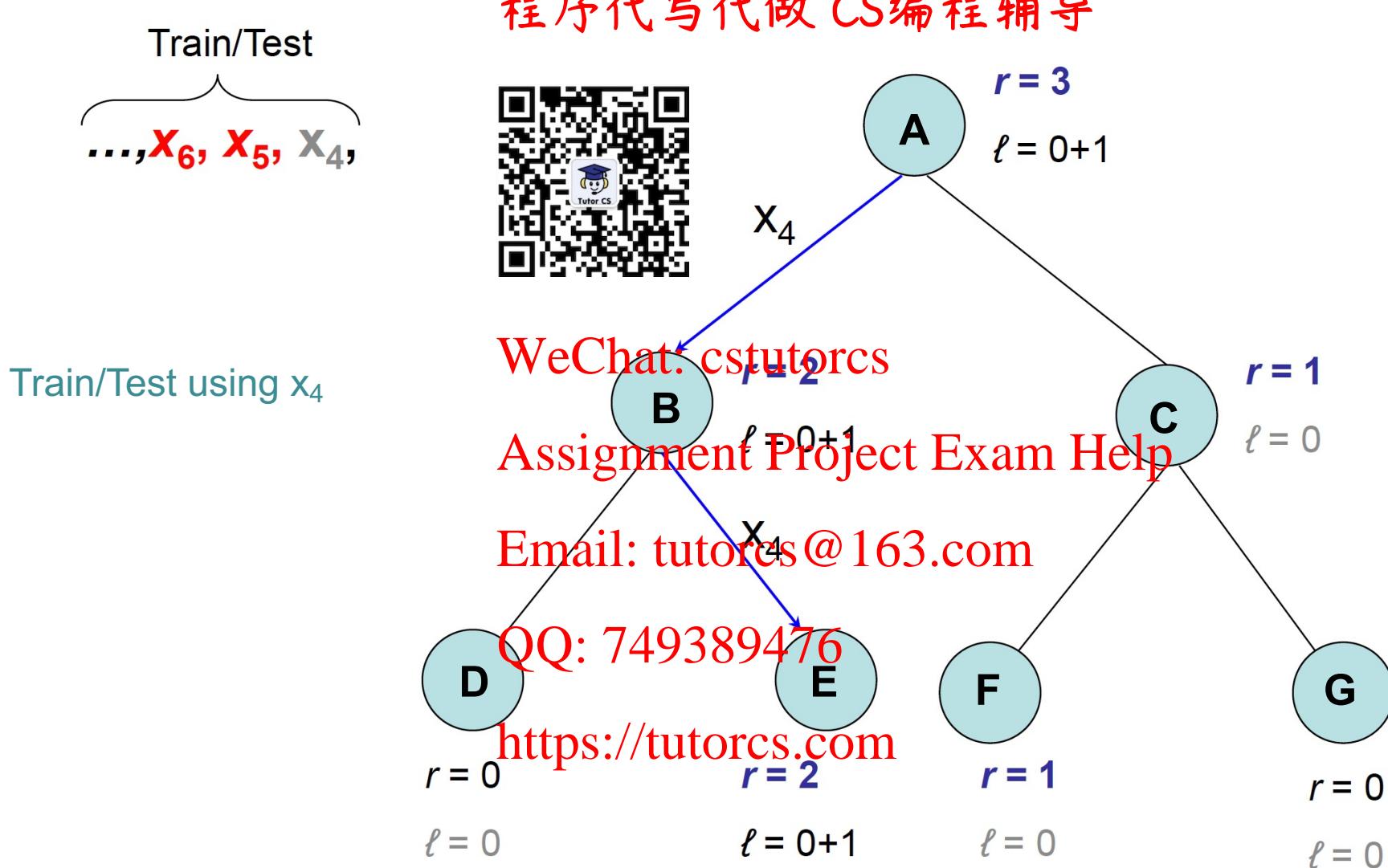
Email: **tutorcs@163.com**

QQ: **749389476** **W2**

https://tutorcs.com
 $X_6, X_5, X_4,$

W1
 X_3, X_2, X_1

Streaming HS-Trees – Example

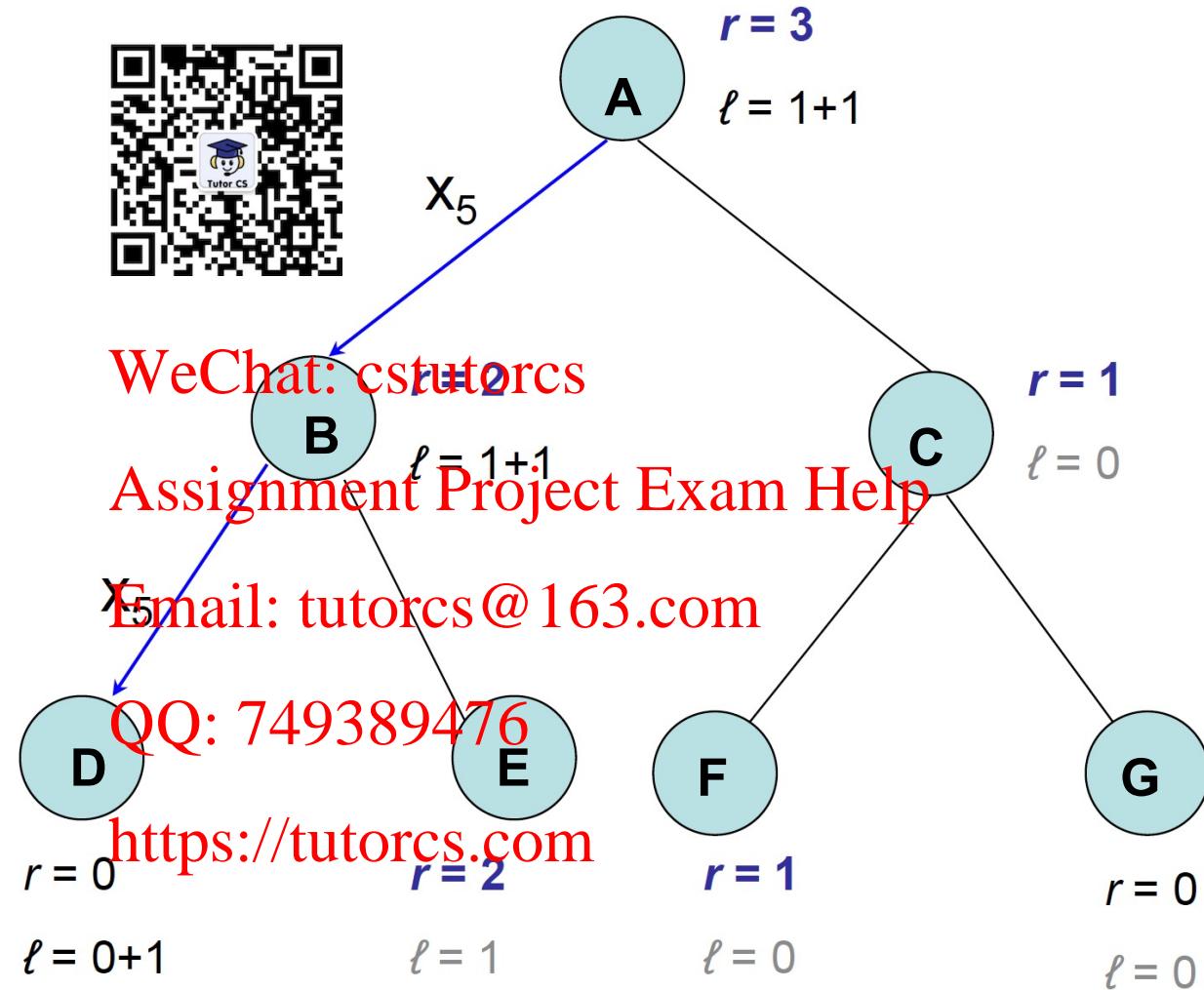


Streaming HS-Trees – Example

Train/Test
 ...
 $x_6, x_5, x_4,$

Train/Test using x_5

程序代写代做 CS 编程辅导



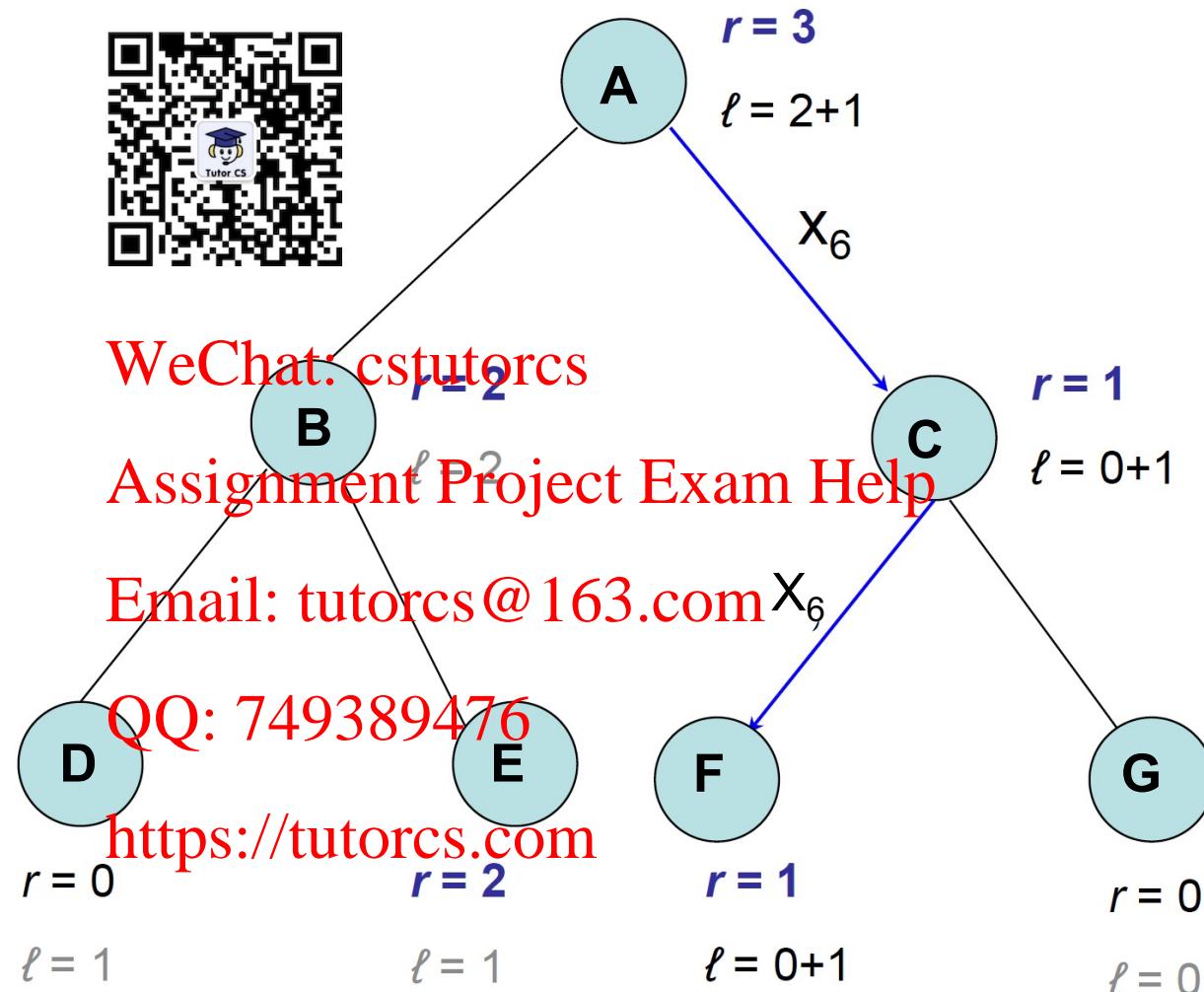
Streaming HS-Trees – Example

程序代写代做 CS 编程辅导

Train/Test

$\dots, x_6, x_5, x_4,$

Train/Test using x_6



程序代写代做 CS编程辅导

When all instances in W2 are processed

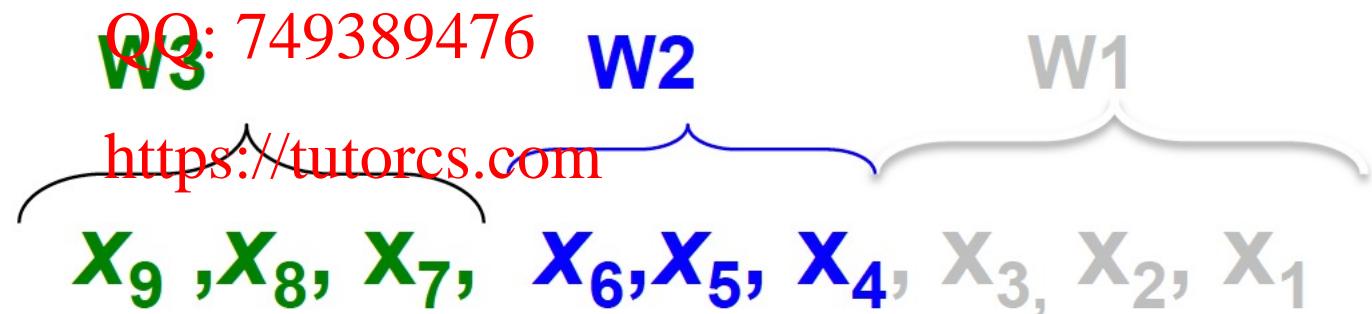
- Model update occurs
- W2 becomes the new reference window
 - Transfer all mass l values to mass r values
 - Reset all mass l values to zero
- W3 becomes the latest window
 - Instances in W3 for training HS-Trees (mass l)
 - Instances in W3 for testing HS-Trees (mass r)

WeChat: cstutorcs

Assignment Project Exam Help

And so on...

Email: tutorcs@163.com



Model Update

程序代写代做 CS编程辅导

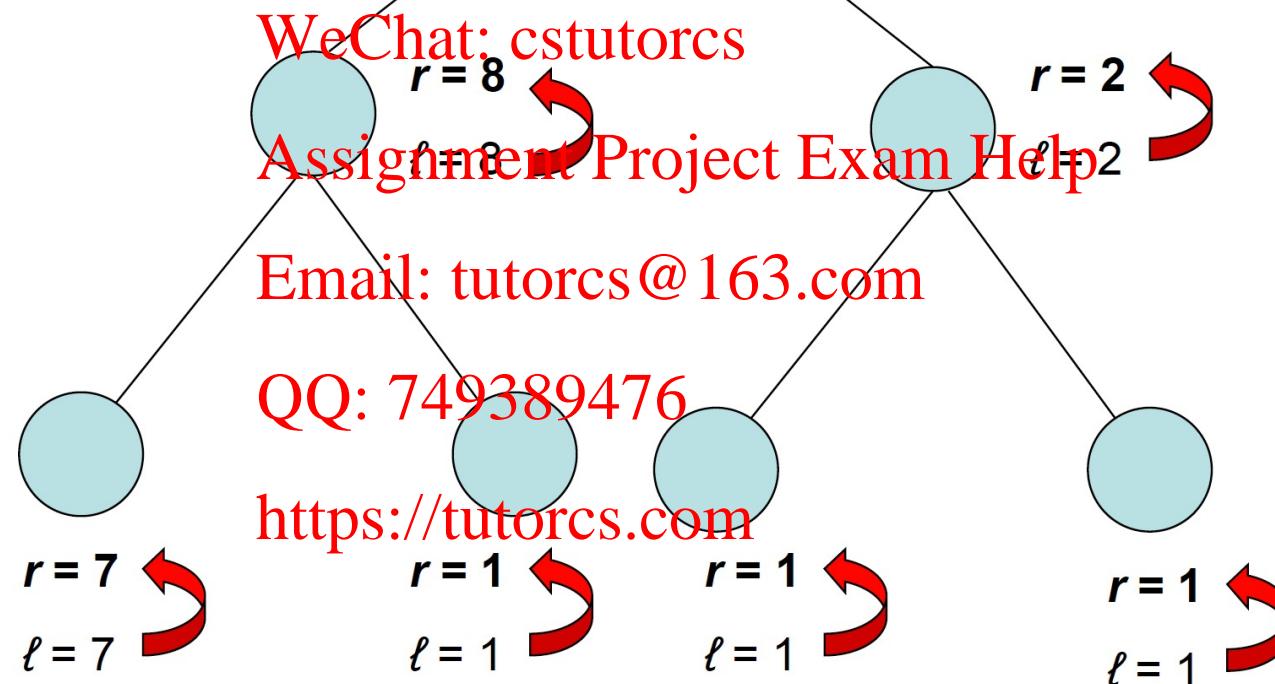
Step 1:

for each node,

$$r \leftarrow \ell$$



$$r = 10 \\ \ell = 10$$



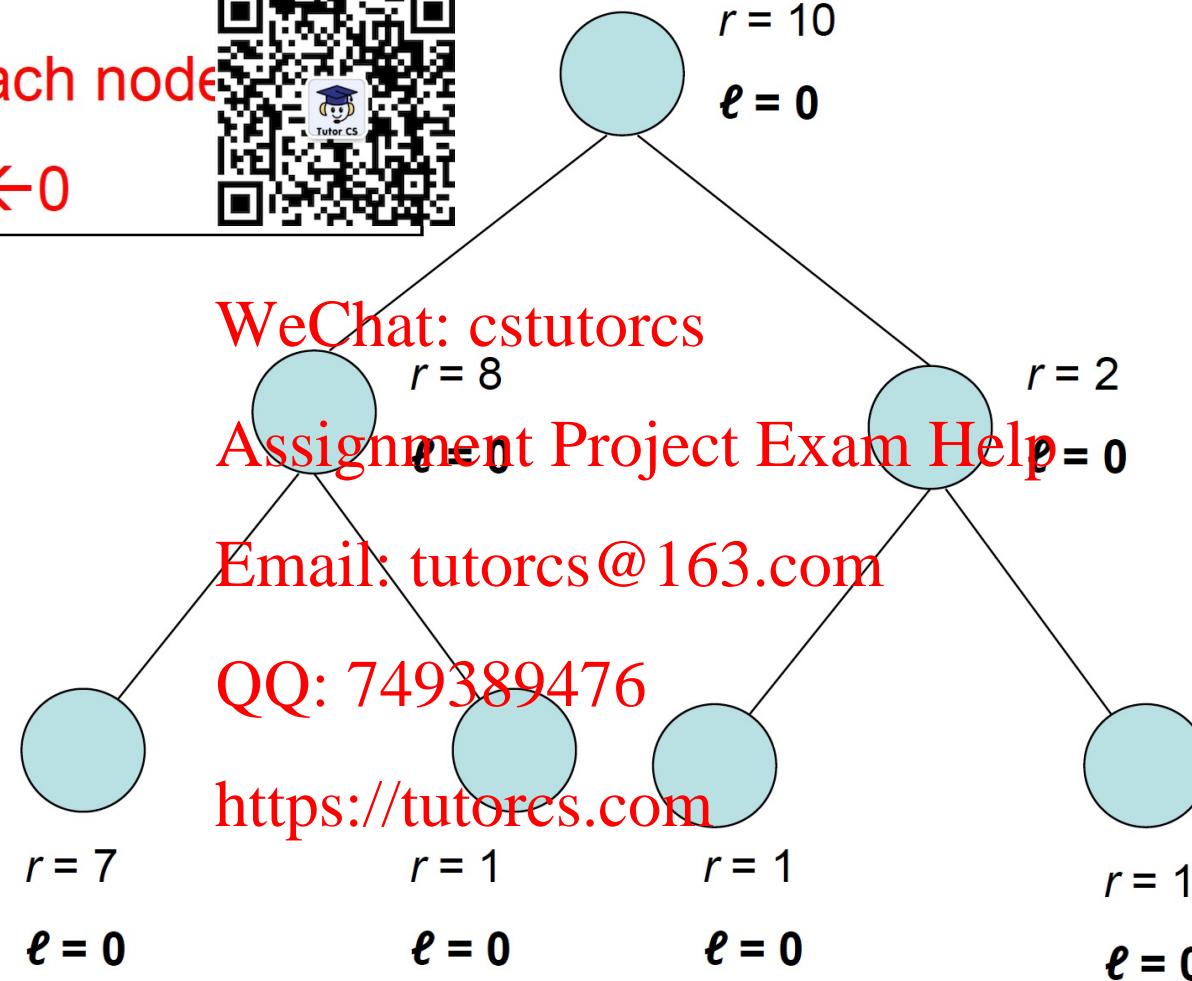
Model Update

程序代写代做 CS编程辅导

Step 2:

for each node

$\ell \leftarrow 0$



Anomaly Score in HS-Tree

- The final score for x is the sum of scores obtained from each HS-Tree in the ensemble



$$\text{anomaly score}(x) = \sum_{t \in T} \text{Score}(x, t_i)$$

WeChat: cstutorcs

$$\text{Score}(x, t_i) = \text{Node}_r \times 2^{\text{Node}_k}$$

Assignment Project Exam Help

Email: tutorcs@163.com
r value at node
QQ: 749389476
Depth of node

<https://tutorcs.com>

程序代写代做 CS编程辅导

Advantages of LOF for anomaly detection:

- Detects anomalies regardless of the data distribution of normal behaviour, since it does not make any assumptions about the distributions of data records.
- Detects anomalies with respect to the density of their neighbouring data records; not to the global model.
- Directly applying LOF to data streams would be extremely computationally inefficient and/or very often may lead to incorrect prediction results.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



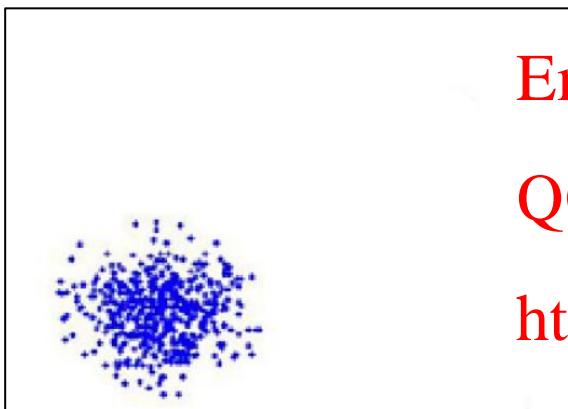
程序代写代做 CS 编程辅导

(i) **Periodic LOF**. Apply LOF algorithm on the entire data set *periodically* (e.g., after every data block of 100² words or 1000 records) or after all the data records are inserted.

- The major problem of this approach is inability to detect anomalies related to the *beginning of new behaviour* that initially appear within the inserted block.

WeChat: cstutorcs

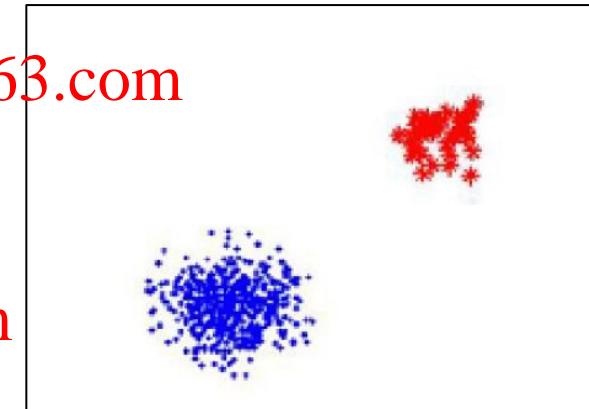
Assignment Project Exam Help



Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



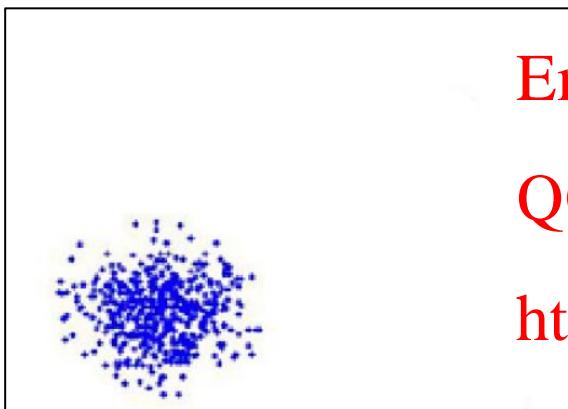
程序代写代做 CS 编程辅导

(ii) **Iterated LOF:** Re-apply the static LOF algorithm *every time a new data record is inserted into the dataset.*

- This static LOF algorithm suffers from the previous problems, but is extremely computationally expensive.
- Increases LOF's time complexity to $O(n^2 \log n)$, where n is the current number of data records in the data set

WeChat: cstutorcs

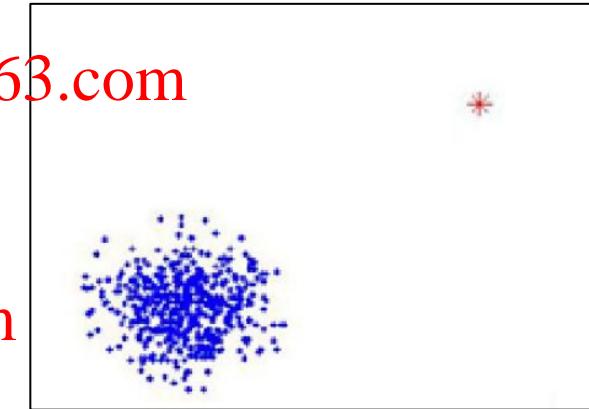
Assignment Project Exam Help



Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

Objectives:

- The result of the incremental algorithm must be equivalent to the result of the “batch”.
- Time complexity of incremental algorithm has to be comparable to the static LOF algorithm $O(n \log n)$.

WeChat: cstutorcs

Step 1 – Insertion:

- **Insertion of new record**, compute $k\text{-dist}$, reachdist , lrd and LOF values of a new point
- **Maintenance**, update $k\text{-dist}$, reachdist , lrd and LOF values for *affected existing points*.

Email: tutorcs@163.com

QQ: 749389476

Step 2 – Deletion: Delete certain data records (e.g., due to their obsoleteness).

- **Maintenance**, update $k\text{-dist}$, reachdist , lrd and LOF values for *affected existing points*.

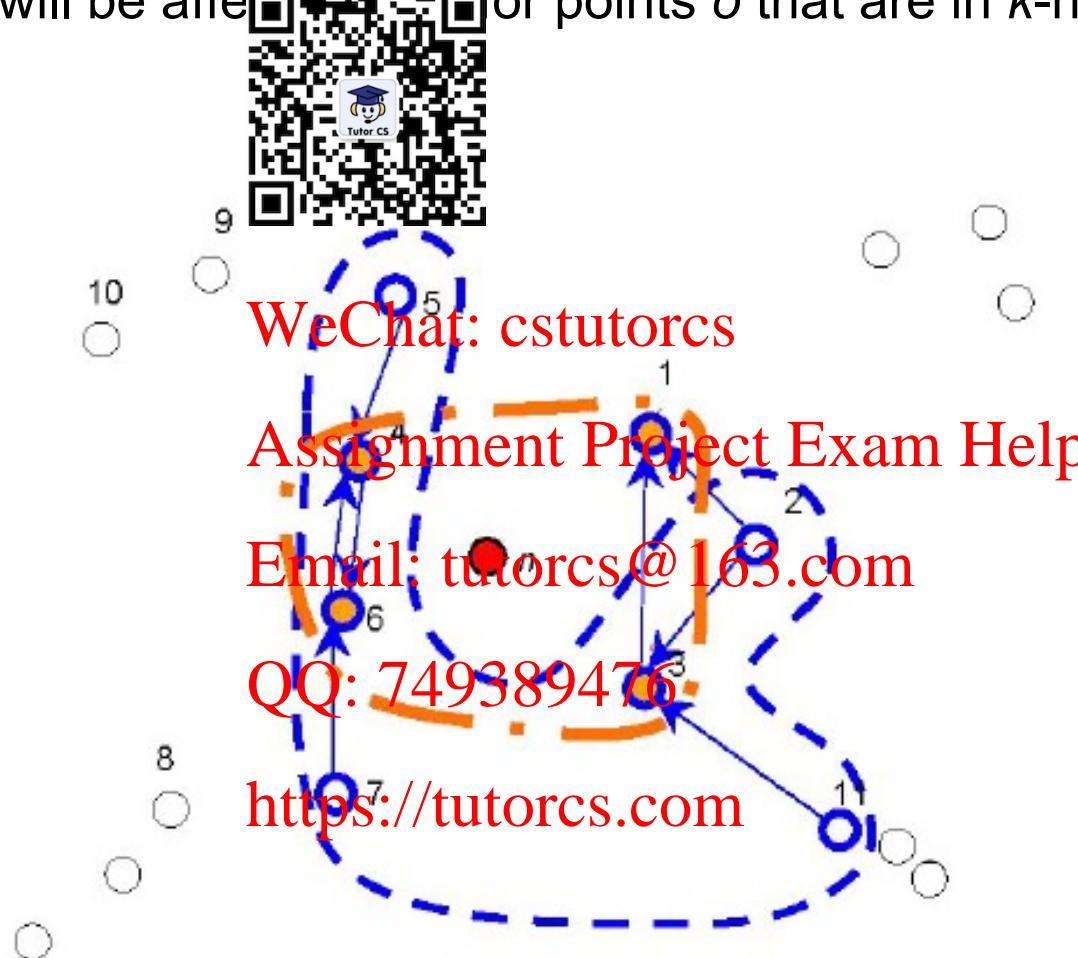
程序代写代做 CS 编程辅导

- **Updating k -dist:** Insertion of the point n may decrease the k -distance of certain neighbouring points, and it is open only to those points that have the new point n in their k -neighbourhood (\dots, k -neighbourhood).
- **Reverse Nearest Neighbour (RNN):** Find all the objects for which the new point n is their (k)-nearest neighbour.



程序代写代做 CS编程辅导

- **Updating $reachdist_k$:** When k -distance(p) changes for a point p , $reachdist_k(o,p)$ will be affected for points o that are in k -neighbourhood of the point p .



程序代写代做 CS 编程辅导

- **Updating lrd :** lrd value of a point p is affected if:
 - The k -neighbourhood of point p changes,
 - $Reachdist$ from point p to each of its k -neighbours changes.



程序代写代做 CS 编程辅导

- **Updating LOF Values:** LOF values of an existing point p should be updated if
 - $lrd(p)$ is updated, or
 - $lrd(p)$ of one of its k -neighbors changes



程序代写代做 CS 编程辅导

- **Updating k -dist:** The deletion of each record p_c from the dataset influences the k -distances of its RNN.
 - k -neighbourhood increases by one each data record p_j that is in reverse k -nearest neighbourhood of p_c . The new k -distance for p_j becomes equal to its distance to its new k -nearest neighbour.

WeChat: cstutorcs

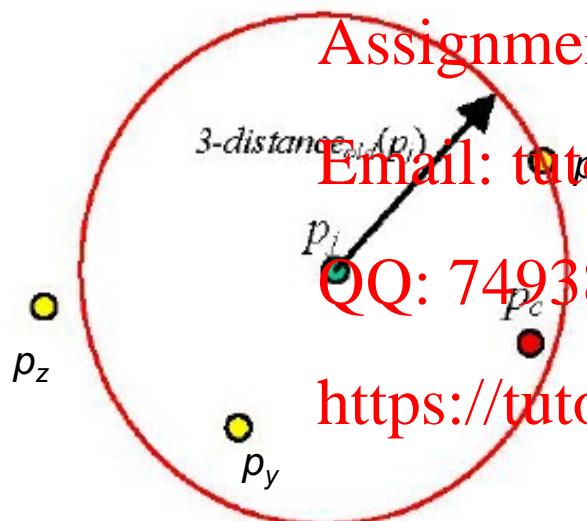
Assignment Project Exam Help

3-distance _{p_{jL}} (p_j)

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS 编程辅导

- **Updating k -dist:** The deletion of each record p_c from the dataset influences the k -distances of its RNN.
 - k -neighbourhood increases by one each data record p_j that is in reverse k -nearest neighbourhood of p_c . The new k -distance for p_j becomes equal to its distance to its new nearest neighbour.

WeChat: cstutorcs



程序代写代做 CS 编程辅导

- **Updating $reachdist$:** The reachability distances from p_j 's nearest neighbours need to be updated.
- **Updating lrd :** lrd value needs to be updated for
 - All points p_j , which k -NN of p_j is updated.
 - All points p_i , which is in k -NN of p_j and p_i is in k -NN of p_j .
- **Updating LOF Values:** LOF value is updated for
 - All points p_j , which lrd value is updated
 - All points p_i , which is in k -NN of p_j



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

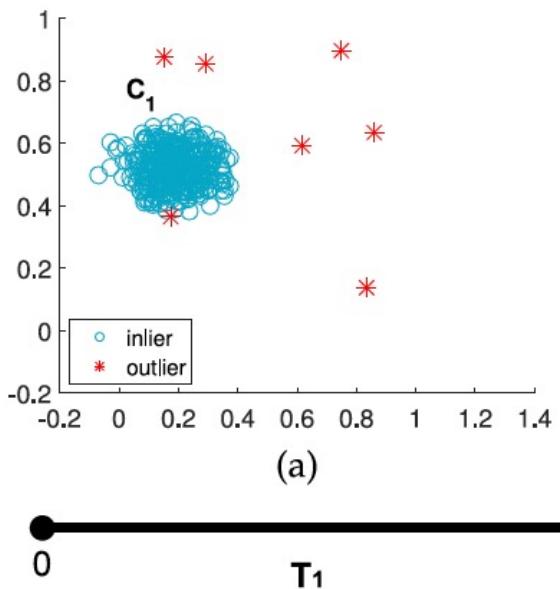
<https://tutorcs.com>

Shortcomings of iLOF

程序代写代做 CS编程辅导

Deleting past data points due to memory limitations causes two problems:

- Differentiation between new events
- The accuracy will drop by one history



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

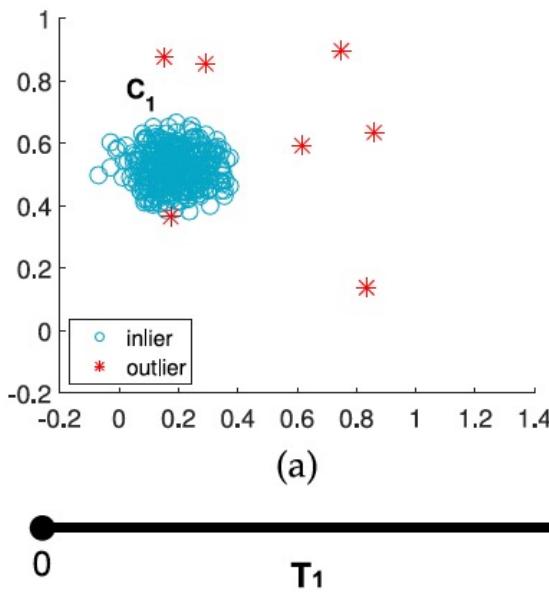
<https://tutorcs.com>

Shortcomings of iLOF

程序代写代做 CS编程辅导

Deleting past data points due to memory limitations causes two problems:

- Differentiation between new events
- The accuracy will drop by one history

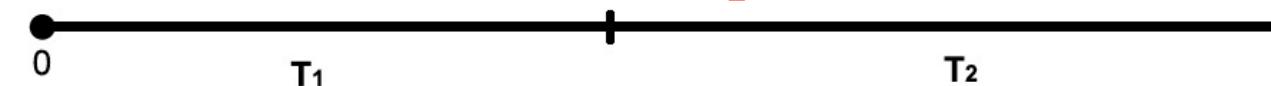
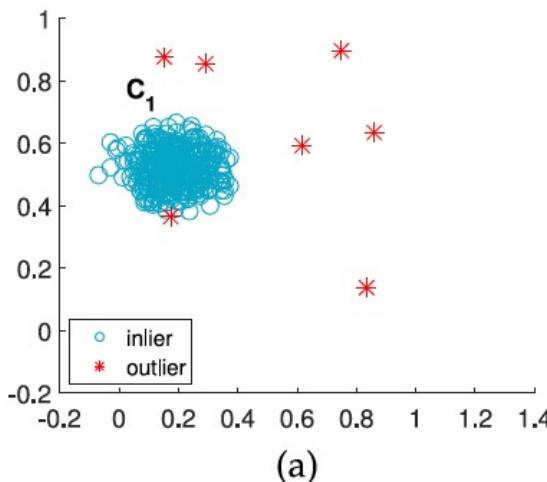


Shortcomings of iLOF

程序代写代做 CS编程辅导

Deleting past data points due to memory limitations causes two problems:

- Differentiation between new events
- The accuracy will drop by one history



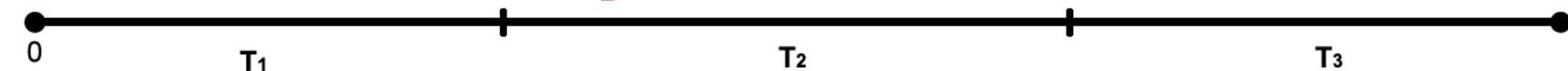
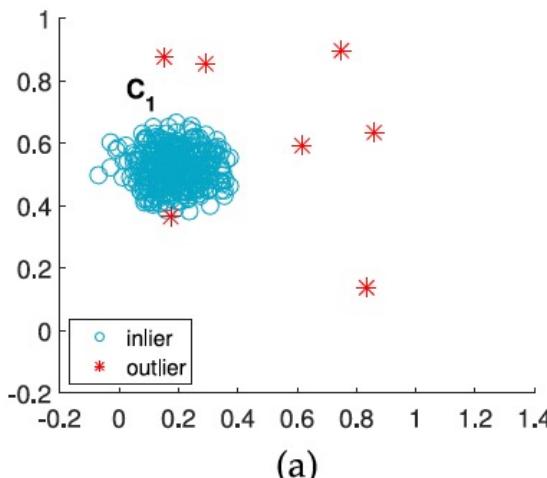
WeChat: cstutorcs
Assignment Project Exam Help
Email: tutorcs@163.com
QQ: 749389476
<https://tutores.com>

Shortcomings of iLOF

程序代写代做 CS编程辅导

Deleting past data points due to memory limitations causes two problems:

- Differentiation between new events
- The accuracy will drop by one history



WeChat: cstutorcs

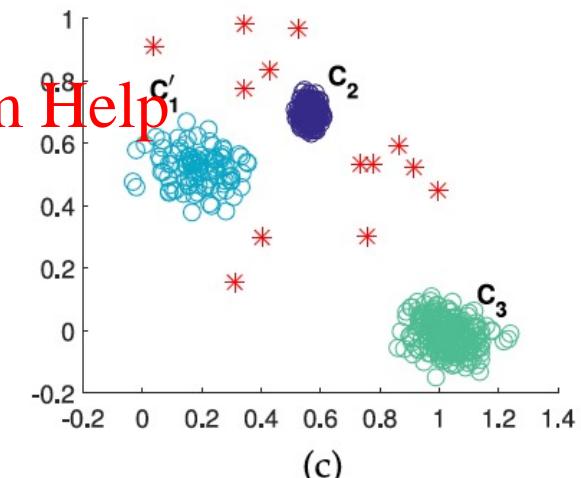
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutores.com>

(b)



程序代写代做 CS编程辅导

Objective: Assign an LOF value to a point p_t , under the constraint that the available memory stores only $m \ll n$ of the n points that have been observed up to time T .

- Need to choose a strategy to summarize the previous data points so that the LOF values of new points can be calculated.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

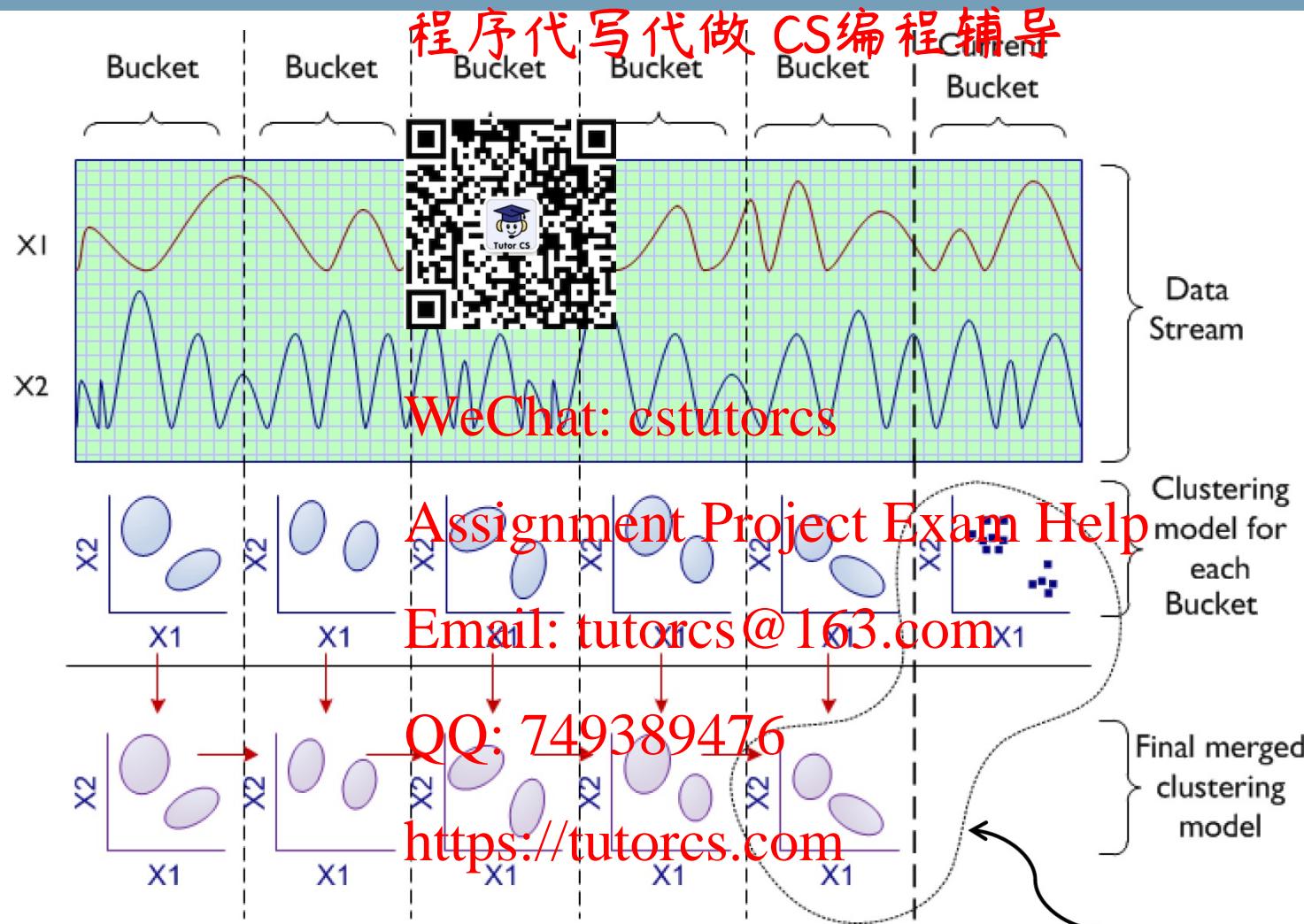
<https://tutorcs.com>



MiLOF Phases:

- Summarization
- Merging
- Revised Insertion

Three Phases of MiLOF – Framework



Compute outlier score of new point using the latest data points and latest merged clustering model

程序代写代做 CS编程辅导

Phase 1 – Summarization:

Build a summary over the past k points along with their corresponding values ($k\text{-}dist$, lrd and LOF), and delete them from memory.



- Every bucket data points are summarized and cluster centres are generated using k -means clustering

WeChat: cstutorcs

Assignment Project Exam Help

Notations:

Email: tutorcs@163.com

- C : points arriving at time T

QQ: 749389476

- Partition C into m clusters $C = \{C_1 \cup C_2 \cup \dots \cup C_m\}$, with cluster centres $V = \{v_1, v_2, \dots, v_m\}$

<https://tutorcs.com>

程序代写代做 CS编程辅导

- **k -dist** of a cluster centre $v_i \in V$



$$k\text{-dist}_k(v_i) = \frac{\sum_{p \in C_i} kdist(p)}{|C_i|}$$

Number of points in C_i

- **lrd** of a cluster centre $v_i \in V$

WeChat: cstutorcs

~~$$lrd_k(v_i) = \frac{\sum_{p \in C_i} lrd_k(p)}{|C_i|}$$~~

Assignment Project Exam Help

Email: tutorcs@163.com

- **LOF** of a cluster centre $v_i \in V$

QQ: 749389476

~~$$lof_k(v_i) = \frac{\sum_{p \in C_i} LOF_k(p)}{|C_i|}$$~~

<https://tutores.com>

程序代写代做 CS编程辅导

Phase 2 – Merging:

Merge the clusters with exist~~ing~~ cluster's to maintain a single set of cluster centres by the anomaly detection framework after each step.

- Using a weighted cluster merging algorithm (weighted k -means) and cluster the cluster centres
- Cluster centre's weight is equal to the *number of data points* in that cluster



WeChat: cstutorcs

Phase 3 – Revised Insertion:

- Compute LOF value of the new incoming data

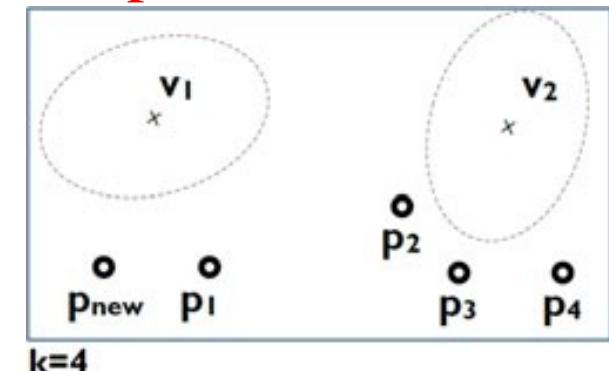
Assignment Project Exam Help

point p , w.r.t. both the recent *data points* and *cluster centres*.

- If a cluster centre is the i^{th} NN of p , we stop searching for the rest of the nearest neighbours.

QQ: 749389476
<https://tutorcs.com>

- Update the *kdist*, *reachdist*, *lrd* and *LOF* values for the existing data points



程序代写代做 CS编程辅导

- What are different windowing techniques for data streams?
- How to apply tree based anomaly detection methods to data streams?
- How to extend LOF for incremental learning while maintaining its performance?



WeChat: cstutorcs

Assignment Project Exam Help

Next: Anomaly Detection Using Support Vector Machine

Email: tutors@163.com

QQ: 749389476

<https://tutorcs.com>

References

程序代写代做 CS编程辅导

1. Swee Chuan Tan, Kai Ming Ting, Tony Fei Liu, "Fast Anomaly Detection for Streaming Data", International Joint Conference on Artificial Intelligence (IJCAI), 2011
 - <https://github.com/yuefei-liu/FastAnomalyDetection>
2. Dragoljub Pokrajac, Aleksandar Lazarevic, Longin Jan Latecki, "Incremental Local Outlier Detection for Data Streams", IEEE Symposium on Computational Intelligence and Data Mining, 2007
3. Mahsa Salehi, Christopher Leckie, James C. Bezdek, Tharshan Vaithianathan, Xuyun Zhang, "Fast Memory Efficient Local Outlier Detection in Data Streams", IEEE Transactions on Knowledge and Data Engineering (TKDE), 2016

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Tutor CS

WeChat: cstutorcs