



程序代写代做 CS编程辅导

An Introduction to 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

- Using machine learning in cybersecurity



- Basics of machine learning

- Introduction to anomaly detection

WeChat: cstutorcs

- Isolation Forest (iForest)

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Intersecting Machine Learning and Cybersecurity



By WIREs Authors



Posted on May 7, 2019

Artificial intelligence in cyber security market is valued at \$4.94bn in 2019, according to Visiongain

GlobeNewswire • May 8, 2019



How AI Beef Up Cybersecurity

Artificial intelligence gives chief information security officers an important new advantage in the ongoing efforts to improve cybersecurity. Find out what to consider when evaluating the latest tools.

程序代写代做 CS编程辅导



Applying AI And Machine Learning To Boost Cybersecurity



Dr. Rao Papolu Forbes Councils
Forbes Technology Council CommunityVoice ⓘ

WeChat: cxtutorcs

Assignment Project Exam Help Automation in Cybersecurity Key to Addressing Growing Risks

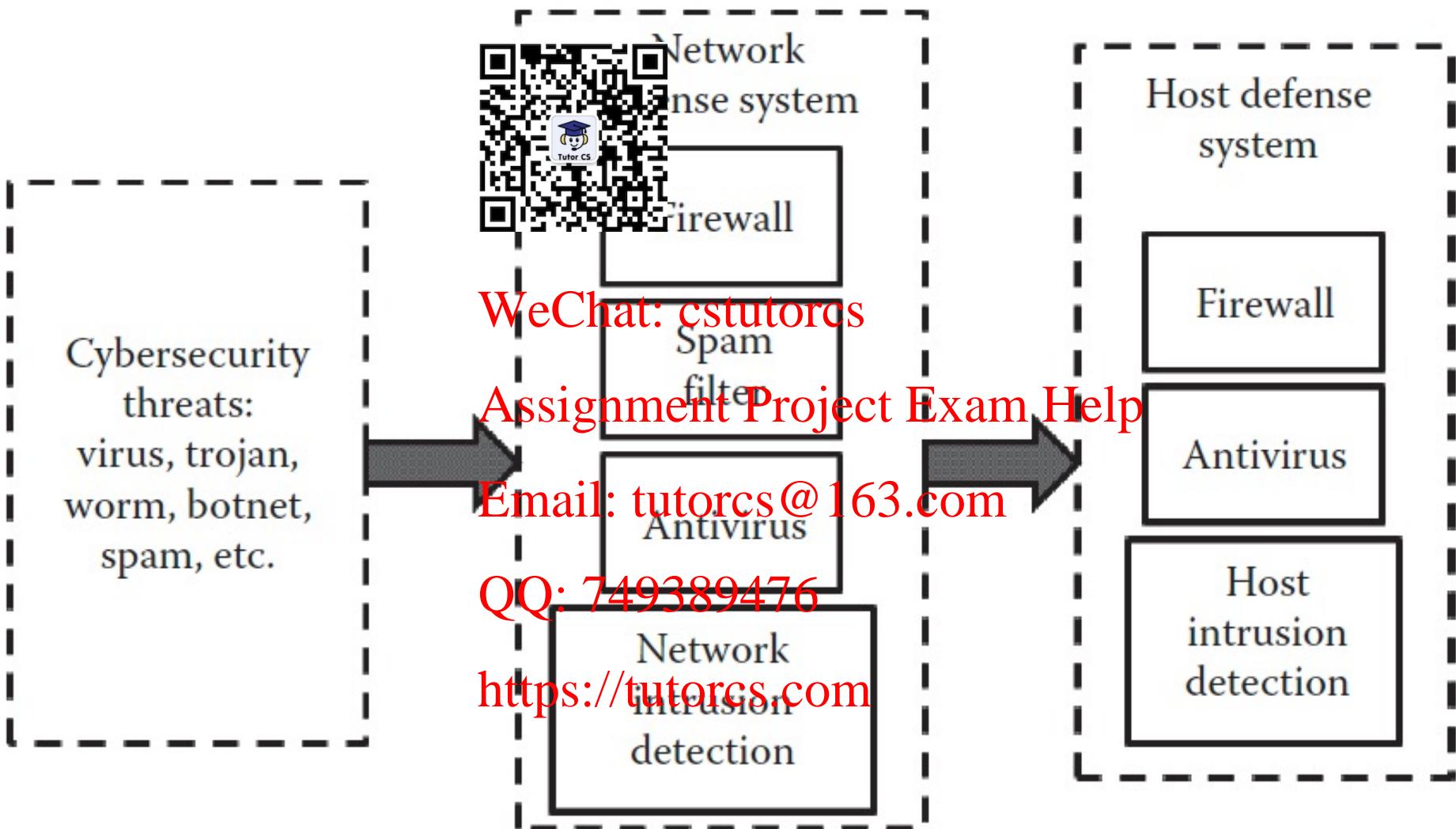
By Simon Eid, Area Vice President, Australia and New Zealand
QQ: 749389476

Simon Eid (CSO Online) on 14 May, 2019 14:29

<https://tutorcs.com>

Conventional Cybersecurity System

程序代写代做 CS编程辅导



Adaptive Defense System for Cybersecurity

程序代写代做 CS编程辅导

Capturing tools

Data preprocessing

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

Analyzing Engine
QQ: 749389476

<https://tutorcs.com>
Decision Process

程序代写代做 CS编程辅导

- **Proactive:**

Maintain the overall security of a system, even if individual components of the system have been compromised by an attack, i.e., *Privacy Preserving Data Mining (PPDM)*.



- **Reactive:**

WeChat: cstutorcs

Identify any unauthorized attempt to access, manipulate, modify, or destroy information or to use a computer system remotely to spam, hack, or modify other computers, i.e., *Intrusion Detection System (IDS)*.

Assignment Project Exam Help
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- **Signature (Misuse) Detection:** Measures the similarity between input events and the signatures of known attacks
- **Anomaly Detection:** Triggers alarms when the detected object behaves significantly differently from predefined normal patterns



Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Overview of ML – Revision

程序代写代做 CS 编程辅导

Input to a machine learning system can consist of instance/measurements about individual entities/objects, e.g. network packet.

- **Attribute** (aka Feature, e.g. variable): component of the instances source IP, destination IP, port, destination port, etc.
- **Label** (aka Response, dependent variable): an outcome that is categorical, numeric, etc. attack vs. legitimate traffic
- **Models**: discovered relationship between attributes and/or label

WeChat: [tutorcs](#)
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- **Supervised learning**

- Teach the computer something (by example), then let it do it
- Use its new-found knowledge to do it
- Labelled data: for given inputs, provide the expected output (“the answer”)
- Infer a function mapping from inputs to outputs



WeChat: cstutorcs

Assignment Project Exam Help

- **Unsupervised learning**

Email: tutorcs@163.com

- Let the computer learn how to do something

QQ: 749389476

- Determine structure and patterns in data

<https://tutorcs.com>

- Unlabelled data: Don't give the computer “the answer”

程序代写代做 CS编程辅导

- **Holdout:** Train a classifier over a fixed training dataset, and evaluate it over a fixed held-out test dataset 
- **Random Subsampling:** Holdout over multiple iterations, randomly selecting the training and testing datasets (maintaining a fixed size for each dataset) on each iteration 
- **Leave-One-Out:** Choose each data point as test case and the rest as training data
- **M-fold Cross-Validation:** Partition the data into M (approximately) equal size partitions, and choose each partition for testing and the remaining M-1 partitions for training

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

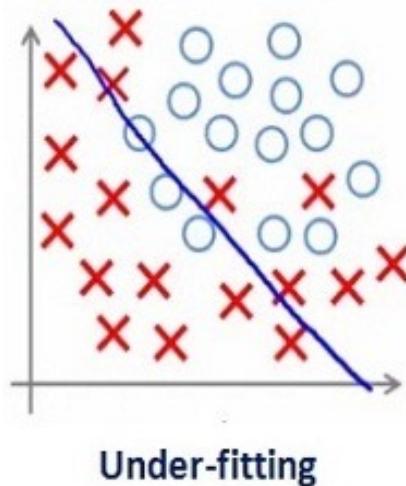
QQ: 749389476

Choose a validation model that is efficient, and minimises bias and variance in evaluation.

<https://tutorcs.com>

Generalisation Problem – Revision

程序代写代做 CS编程辅导



WeChat: cstutorcs

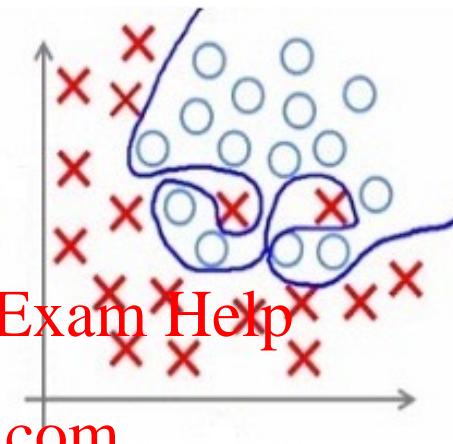
Assignment Project Exam Help

Email: tutorcs@163.com

Appropriate-fitting

QQ: 749389476

<https://tutorcs.com>



- Confusion Matrix



		Actual	
		Cat	Dog
Predicted	Cat	4 (TP)	3 (FP)
	Dog	2 (FN)	6 (TN)

Assignment Project Exam Help

Email: tutorcs@163.com

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \text{ QQ: 749389476}$$

$$= \frac{4 + 6}{4 + 6 + 3 + 2} \cong 67\%$$

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Anomaly detection

- Number of negative examples = 9990
 - Number of positive examples = 10



- If model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$

- Accuracy is misleading because model does not detect any positive examples

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

- **Recall:**

$$\frac{TP}{TP + FN}$$



- **Precision:**

$$\frac{TP}{TP + FP}$$

WeChat: cstutorcs

Assignment Project Exam Help

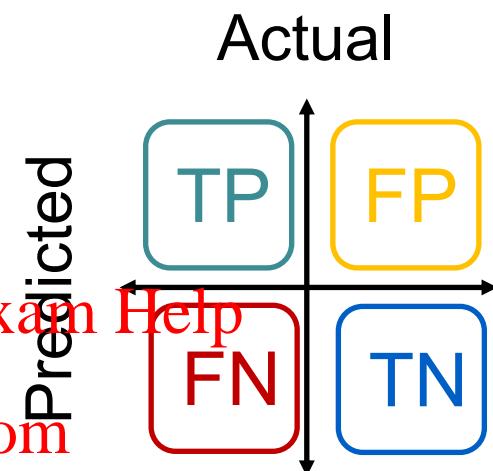
Email: tutorcs@163.com

- **F-Score:**

$$(1 + \beta^2) \frac{Per \times Rec}{Rec + \beta^2 Per}$$

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

程序代写代做 CS编程辅导



ROC (Receiver Operating Characteristic) Curve

- ROC curve plots TPR (on the y-axis) against FPR (on the x-axis)
- Performance of each classifier represented as a point on the ROC curve
- (TPR,FPR):
 - (0,0): declare everything to be negative class
 - (1,1): declare everything to be positive class
 - (1,0): ideal
- Diagonal line:
 - Random guessing
 - Below diagonal line:
 - prediction is opposite of the true class



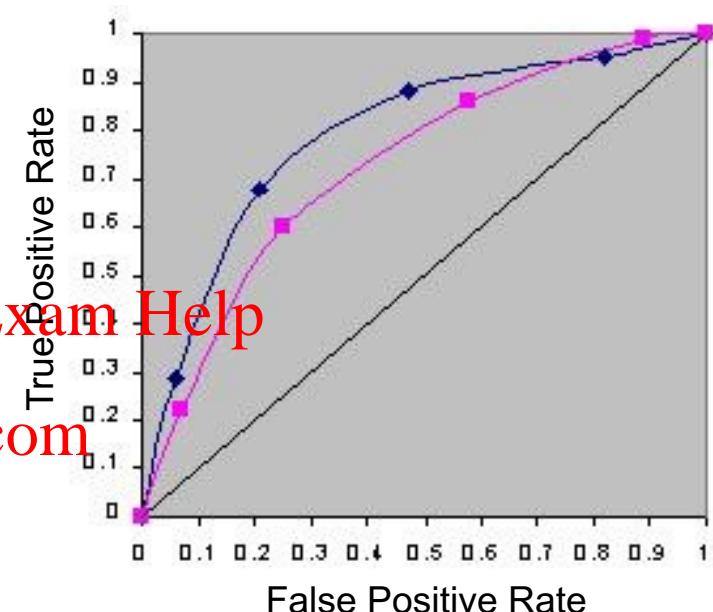
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Anomaly (Outlier) Detection

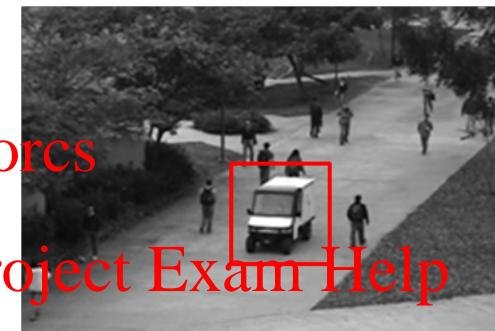
程序代写代做 CS编程辅导

- An anomaly is defined as a pattern in data that does *not conform to the expected behaviours*, including outliers, abbreviations, contaminants, and surprise, etc., in applications.



No.	Delta Time	Source	Destination	Protocol
1	0.000000	FC:AA	FF:FF:FF	ARP Request
2	0.028674	00:1C	FF:FF:FF	ARP Request
3	0.200358	FC:AA	33:33:00	LLMNR
4	0.000001	192.168.1.10	224.0.0.1	LLMNR
5	0.083002	192.168.1.10	234.12.1.1	LLMNR
6	0.016290	FC:AA	33:33:00	LLMNR
7	0.000001	192.168.1.10	224.0.0.1	LLMNR

WxChat: cstutorcs
Assignment Project Exam Help



Email: tutorcs@163.com



Types of Anomalies

程序代写代做 CS 编程辅导

- **Global (Point) Anomalies:** A data object is a global outlier if it *deviates significantly from the rest of the data set*. To detect global anomalies, the critical issue is to find an appropriate measure of deviation with respect to the application in question.
- **Contextual (Conditional) Anomalies:** A data object is a contextual anomaly if it *deviates significantly with respect to a specific context of the object*. In contextual anomaly detection, the context has to be specified as part of the problem definition.
- **Collective Anomalies:** A subset of data objects forms a collective anomaly if the objects as a *whole deviate significantly from the entire data set*. Importantly, the individual data objects may not be anomalies.



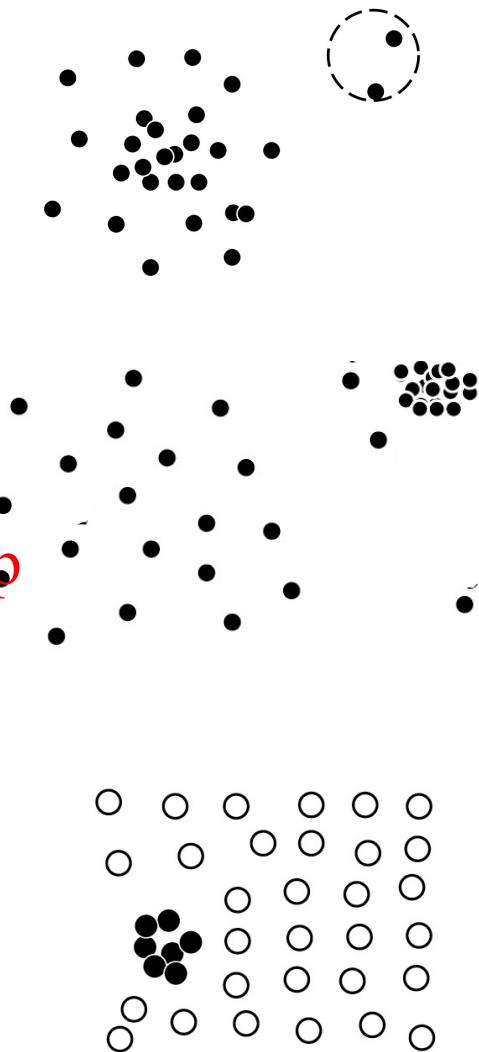
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

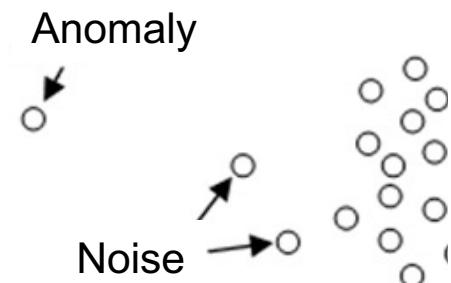
<https://tutorcs.com>



程序代写代做 CS编程辅导

- **Noise vs. Anomaly:**

- Noise is a random error in an instance variable.
- In general, noise is not interesting in data analysis, including anomaly detection.
- Anomalies are interesting because they are suspected of not being generated by the same mechanisms as the rest of the data.



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

- **Noise vs. Anomaly:**

- Noise is a random error in an instance variable.
- In general, noise is not interesting in data analysis, including anomaly detection.
- Anomalies are interesting because they are suspected of not being generated by the same mechanisms as the rest of the data.



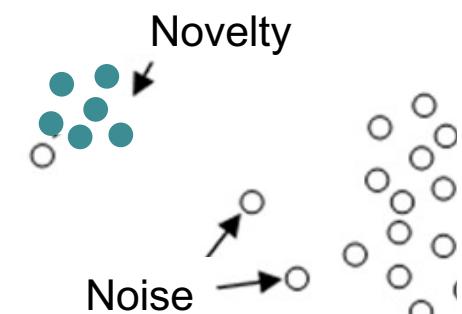
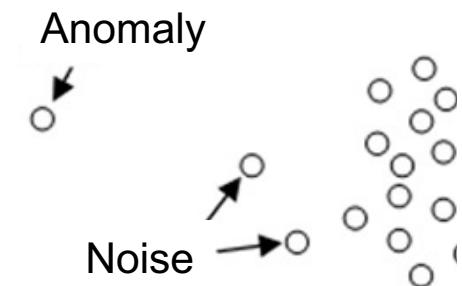
WeChat: cstutorcs

Assignment Project Exam Help

- **Novelty vs. Anomaly:**

Email: tutorcs@163.com

- In evolving datasets, novel patterns may initially appear as anomalies
- Once new patterns are confirmed, they are usually incorporated into the model of normal behaviour so that follow-up instances are not treated as anomalies anymore.



程序代写代做 CS编程辅导

General Steps

- Build a profile of the “normal” population
 - Profile can be patterned from summary statistics for the overall population
- Use the “normal” profile to detect anomalies
 - Anomalies are observations whose characteristics differ significantly from the normal profile



WeChat: cstutorcs

Methods

1. Extreme Value Analysis Email: tutorcs@163.com
2. Proximity-Based
3. Model-based

Assignment Project Exam Help

QQ: 749389476

<https://tutorcs.com>

1. Extreme Value Analysis

程序代写代做 CS编程辅导

- Assume a parametric model describing the distribution of the data (e.g., normal distribution)
- Apply a statistical test (e.g., $(x - \mu)/\sigma$) that depends on
 - Data distribution
 - Parameter of distribution (e.g., mean, variance)
 - Number of expected anomalies (confidence limit)



WeChat: cstutorcs

Limitations

- Most of the tests are for a single attribute
- In many cases, data distribution may not be known
- For high dimensional data, it may be difficult to estimate the true distribution
 - Can be used as final steps for interpreting outputs of other anomaly detection methods

Assignment Project Exam Help

Email: tutors@163.com

QQ: 749389476

<https://tutors.com>

2. Proximity-Based

程序代写代做 CS编程辅导

- Data is represented as a vector of features.
- Assumes the proximity of an object to its neighbourhood significantly deviates from the proximity of the object to most of the other objects in the dataset.
- Three major approaches

WeChat: cstutorcs

Assignment Project Exam Help

2.1 Nearest-neighbour based

Email: tutorcs@163.com

2.2 Density based

QQ: 749389476

2.3 Clustering based

<https://tutorcs.com>



2.1 Nearest-Neighbour Based

程序代写代做 CS 编程辅导

- Compute the distance between every pair of data points
- There are various ways to detect anomalies:



- Data points for which there are fewer than k neighbouring points within a distance D
- The top n data points whose distance to the k^{th} nearest neighbour is greatest
- The top n data points whose average distance to the k^{th} nearest neighbours is greatest

WeChat: cstutors
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

2.2 Density-based

程序代写代做 CS编程辅导

- Estimates the density of objects (using proximity measures between objects).
- Objects that are in regions of low density are relative distance from their nearest neighbours, and can be considered anomalous.
- A more sophisticated approach accommodates the fact that data sets can have regions of widely differing densities.
 - Classifies a point as an outlier only if it has a local density significantly less than that of most of its neighbours.



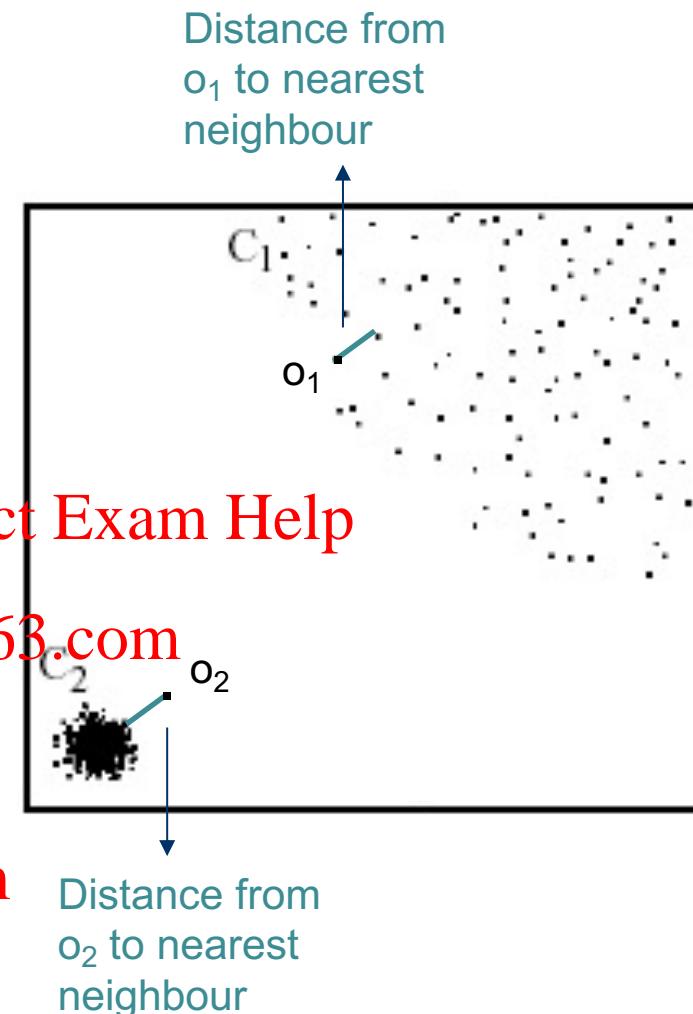
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



2.3 Clustering-Based

程序代写代做 CS编程辅导

- Cluster the data into groups of different density
- Choose points in small clusters candidate anomalies
- Compute the distance between candidate points and non-candidate clusters.



WeChat: cstutorcs

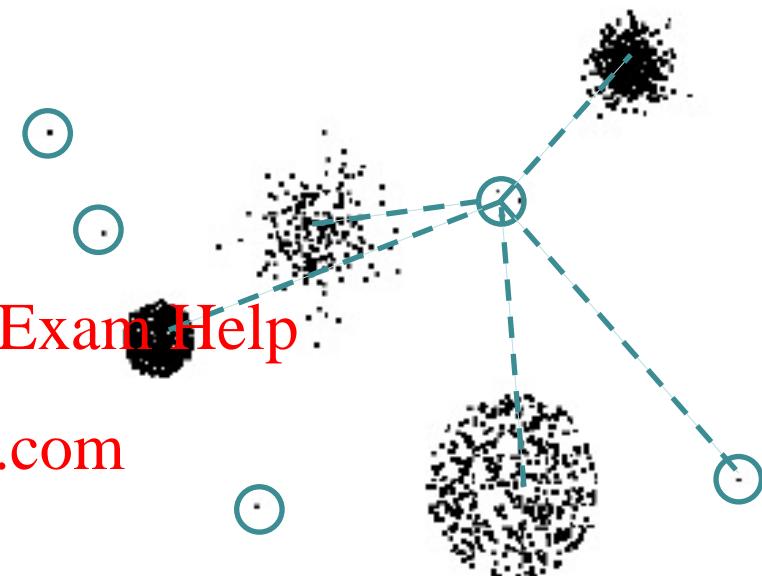
Assignment Project Exam Help

Email: tutorcs@163.com

- If candidate points are far from all other non-candidate points, they are anomalies

QQ: 749389476

<https://tutorcs.com>



3.1 Classification-Based Methods

程序代写代做 CS 编程辅导

Idea: Train a classification model that can distinguish “normal” data from anomalies

- Consider a training set that samples **labelled** as “normal” and others **labelled** as “anomaly”
 - But, the training set is typically heavily biased: number of “normal” samples likely far exceeds number of anomaly samples
- Handle the imbalanced distribution
 - Oversampling positives and/or under sampling negatives
 - Cost-sensitive learning

Assignment Project Exam Help

Email: tutors@163.com

QQ: 749389476

<https://tutorcs.com>



WeChat: cstutorcs

程序代写代做 CS 编程辅导

- One-class model: A classifier is built to describe only the normal class

- Learn the decision boundary for the normal class using classification methods such as one-class SVM
 - Any samples that do not belong to the normal class (not within the decision boundary) are declared as anomalies
 - Advantage: can detect new anomalies that may not appear close to any anomalous objects in the training 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

- **Scoring Techniques:**

- Assign an anomaly score to each instance in the test data depending on the degree to which the instance is considered an anomaly.
- The output is a ranked list of anomalies.
- An analyst may choose to either analyse top few anomalies or use a cut-off threshold (or domain specific threshold) to select the anomalies.

WeChat: cstutorcs
Assignment Project Exam Help

- **Labelling Techniques:** Email: tutorcs@163.com

- Assign a label (normal or anomalous) to each test instance.
QQ: 749389476
- Limit the analysts to the binary label, (though this can be controlled indirectly through parameter choices within each technique).

程序代写代做 CS 编程辅导

- Modelling data with skewed class distributions (class imbalance)
- Sheer volume and heterogeneity of network data
- Difficult to assess the performance of the system, given the vast possibilities of anomalies and lack of labels
- Cost of error in IDS is huge
- Large false alarm rate degrades confidence in the system
- Lack of interpretability
- Anomalies may be undetectable at one level of granularity or abstraction but easy to detect at another level
- Evolving patterns (concept drift)



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



Isolation Forest (iForest) [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

- **Objective:** Isolates anomalies rather than profiles normal instances
- **Isolation:** Separating an  from the rest of the instances

To achieve this, we take advantage of two anomalies' quantitative properties:

- i. They are the minority consisting of fewer instances, and 
- ii. They have attribute-values that are very different from those of normal instances 

Email: tutorcs@163.com

- **Isolation Tree (iTTree) Intuition:** Because of their susceptibility to isolation, anomalies are isolated closer to the root of the tree; whereas normal points are isolated at the deeper end of the tree.

 <https://tutorcs.com>

程序代写代做 CS编程辅导

- Anomalies are more susceptible to isolation under random partitioning

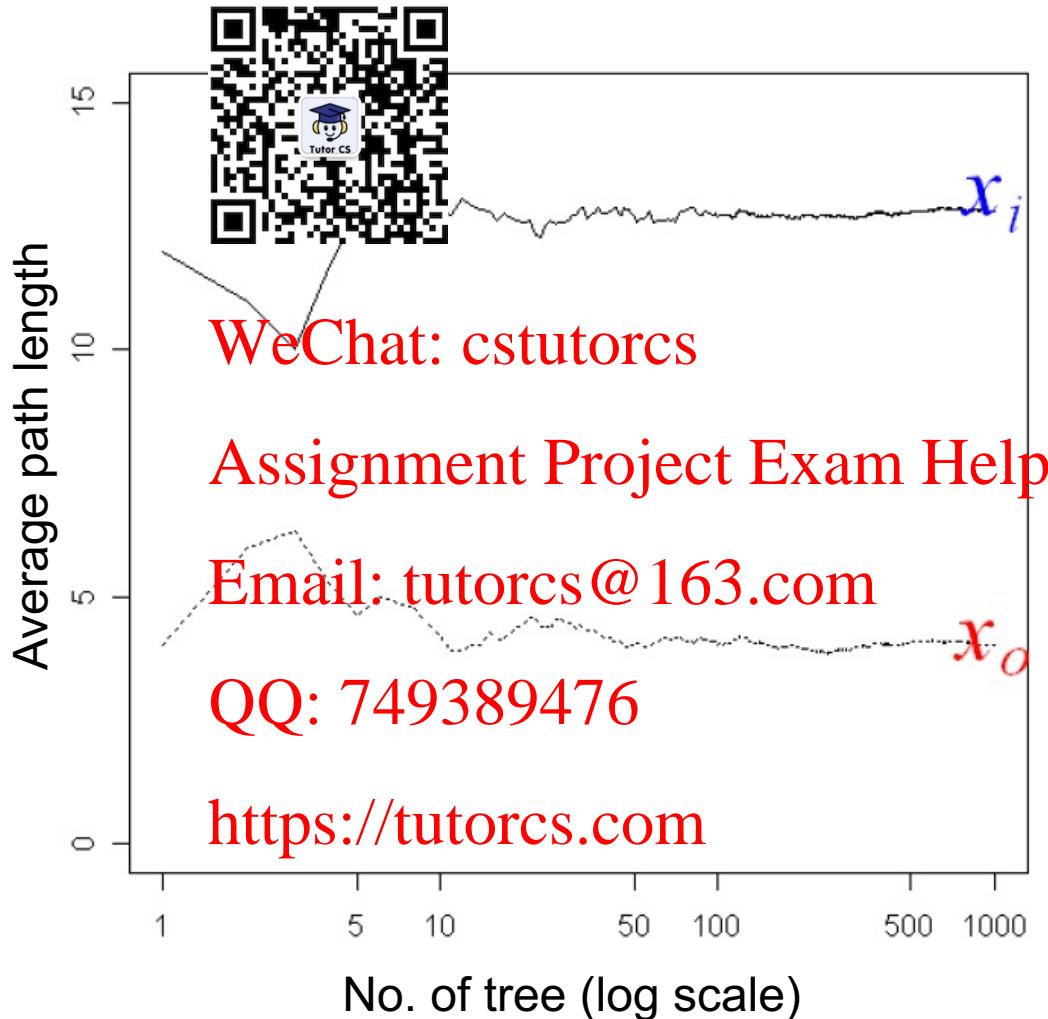


(a) x_i requires only 12 partitions. (b) x_o requires only 4 partitions
<https://tutorcs.com>

Figure. Identifying normal vs. abnormal observations

程序代写代做 CS 编程辅导

- Anomalies are more susceptible to isolation and hence have short path lengths



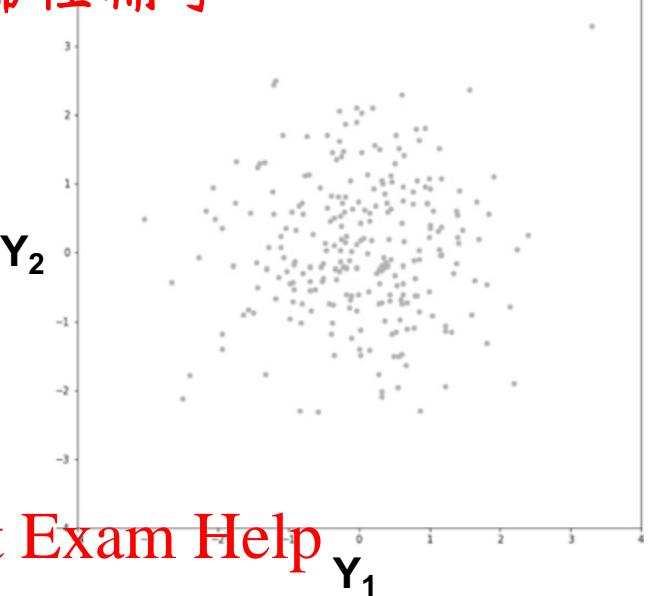
程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a di



WeChat: cstutorcs

Assignment Project Exam Help



Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Isolation Forest

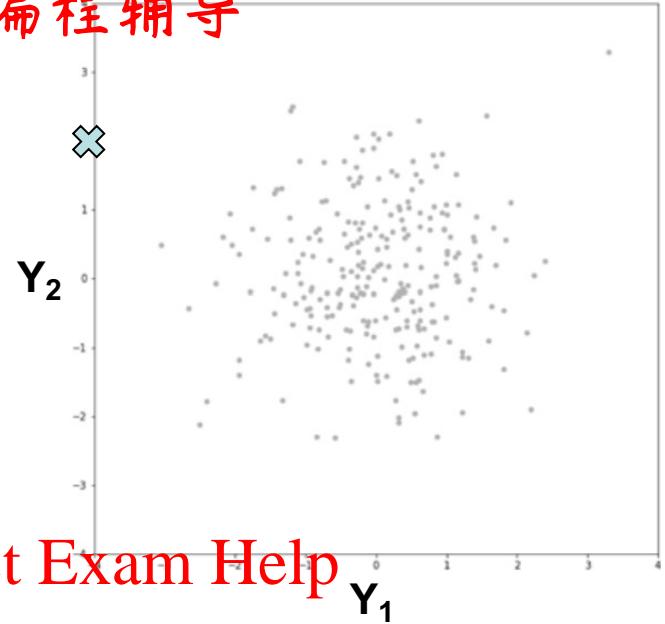
程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value for that dimension



WeChat: cstutorcs

Assignment Project Exam Help



Email: tutorcs@163.com

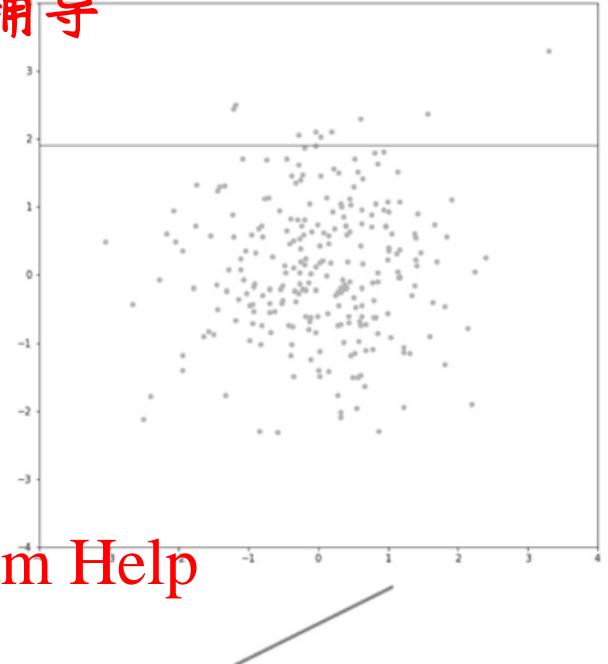
QQ: 749389476

<https://tutorcs.com>

Isolation Forest

程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value
 - Draw a straight line through the data at that value and split data



WeChat: cstutorcs
Assignment Project Exam Help

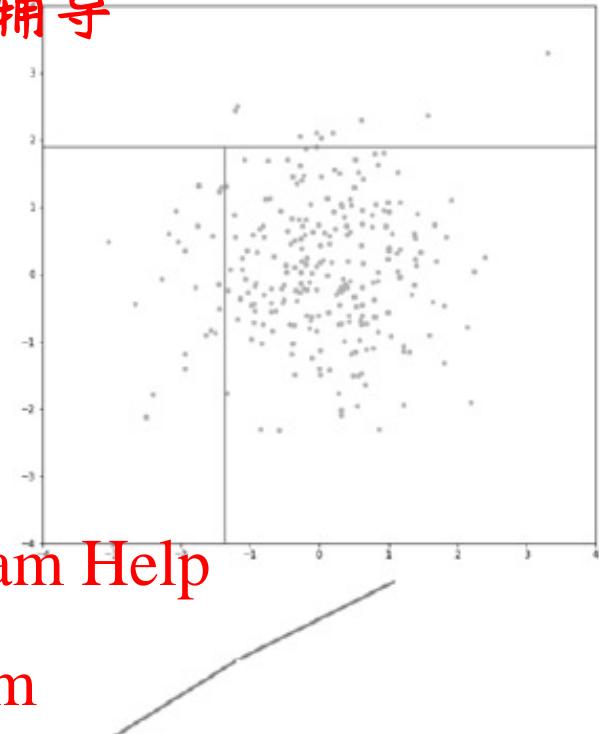
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value
 - Draw a straight line through the data at that value and split data
 - Repeat until tree is complete



Email: tutorcs@163.com

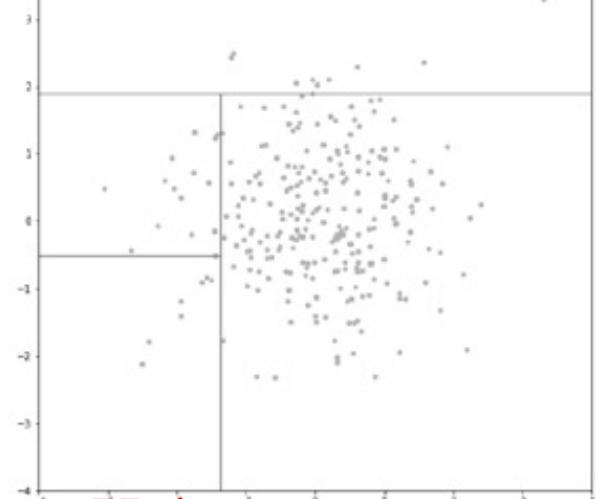
QQ: 749389476

<https://tutorcs.com>

Isolation Forest

程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value
 - Draw a straight line through the data at that value and split data
 - Repeat until tree is complete

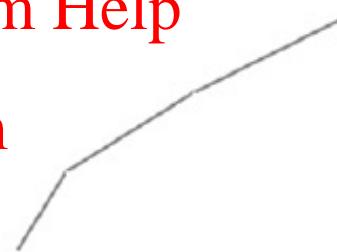


WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

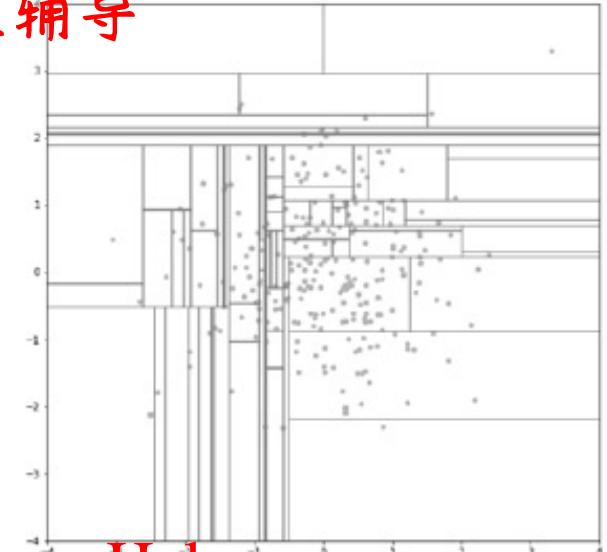
<https://tutorcs.com>



Isolation Forest

程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value from that dimension
 - Draw a straight line through the data at that value and split data
 - Repeat until tree is complete

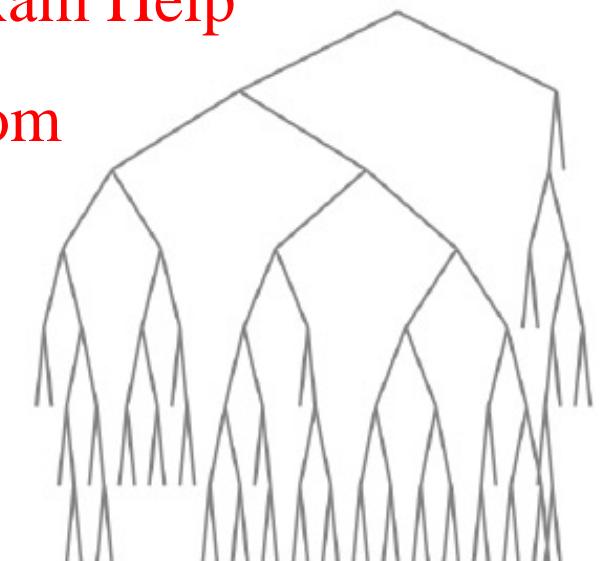


WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Isolation Forest

程序代写代做 CS编程辅导

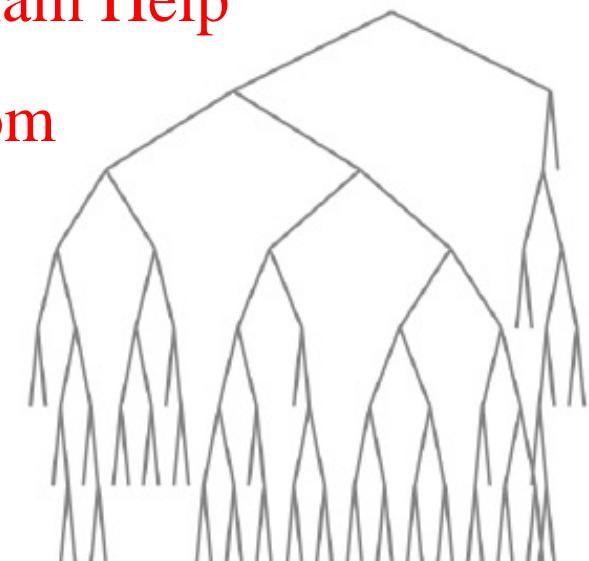
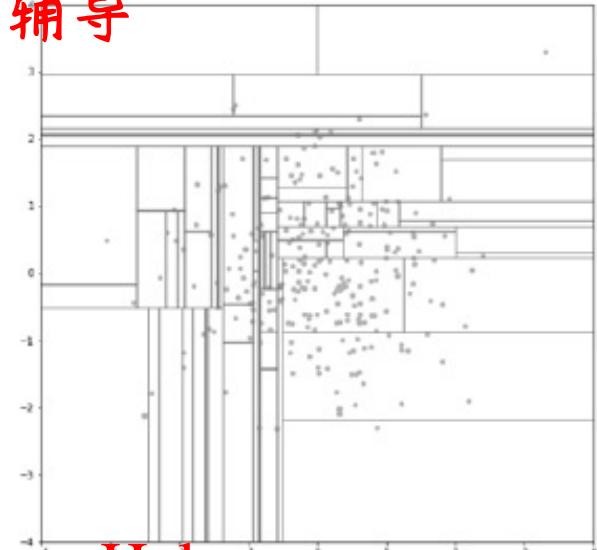
- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value from that dimension
 - Draw a straight line through the data at that value and split data
 - Repeat until tree is complete
- Generate multiple trees → forest

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

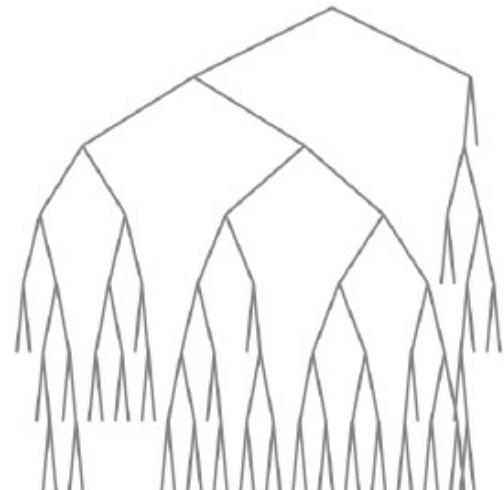
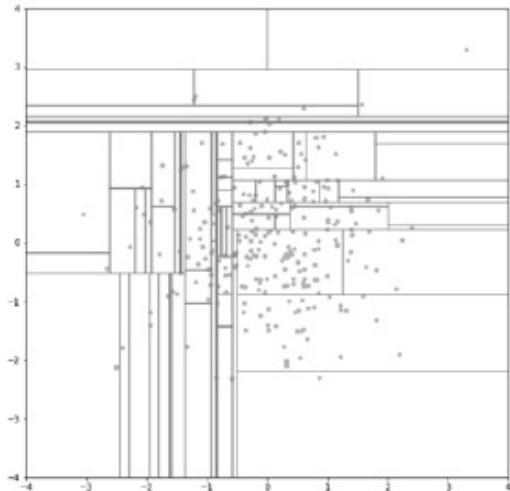
<https://tutorcs.com>



Isolation Forest

程序代写代做 CS编程辅导

iTree 1



iTree 2



WeChat: cstutorcs

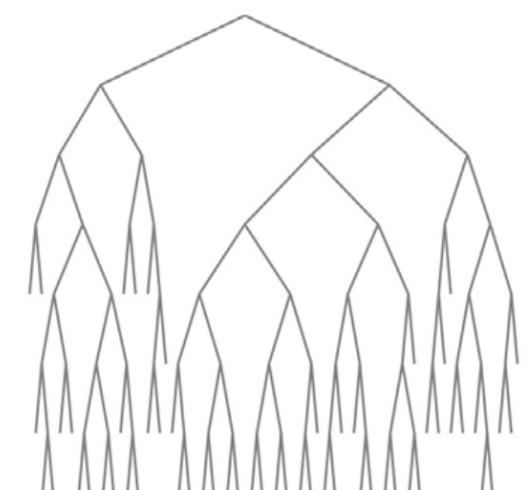
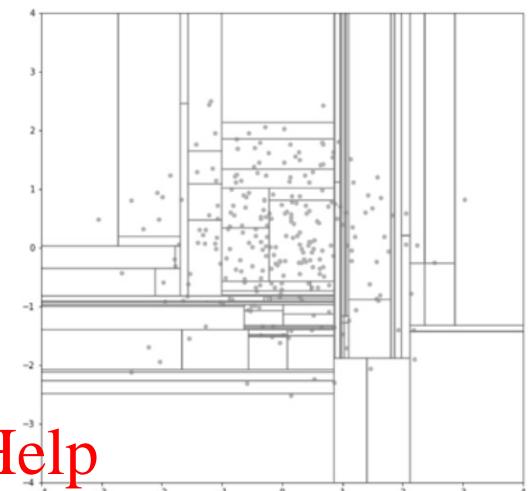
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

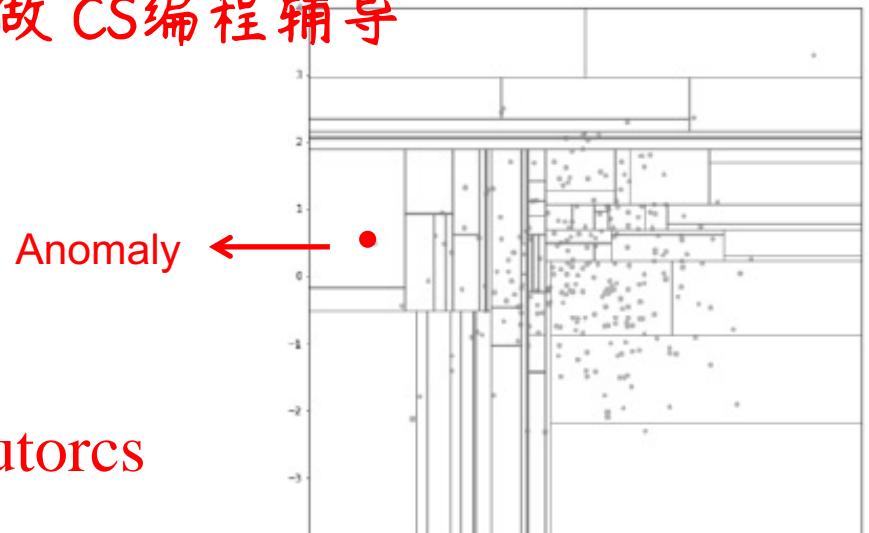
iTree 3



Isolation Forest

程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value
- Draw a straight line through the data at that value and split data
- Repeat until tree is complete
- Generate multiple trees → forest
- Anomalies will be isolated in only a few steps

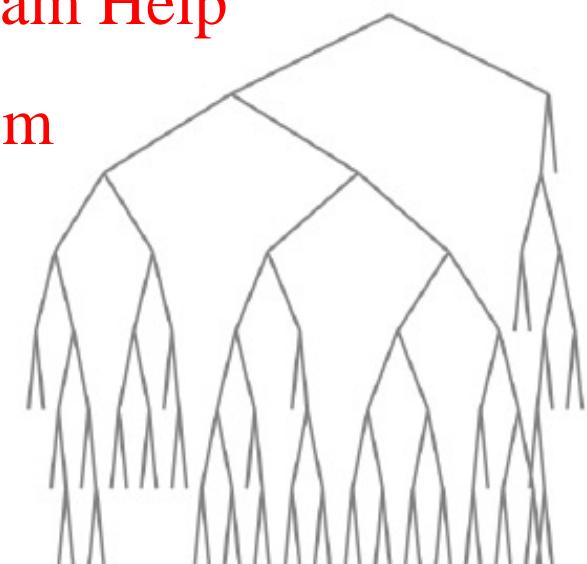


WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

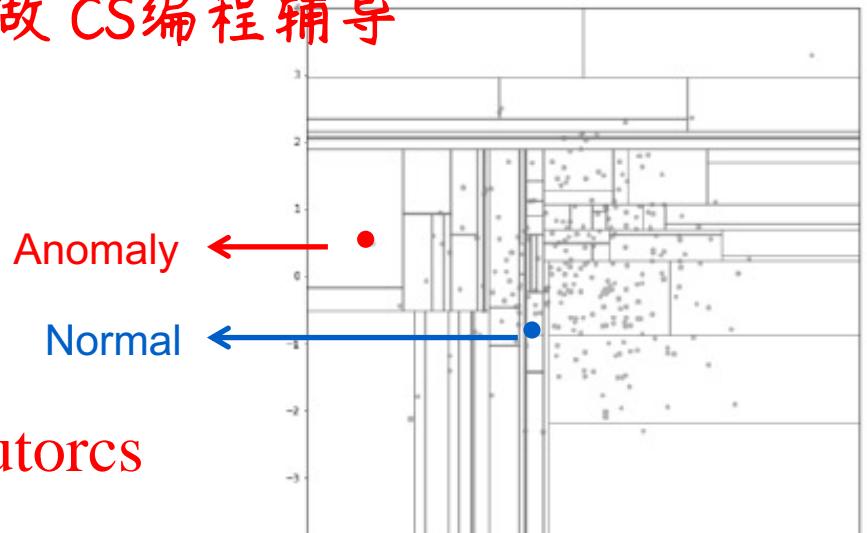
<https://tutorcs.com>



Isolation Forest

程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value
- Draw a straight line through the data at that value and split data
- Repeat until tree is complete
- Generate multiple trees → forest
- Anomalies will be isolated in only a few steps
- Nominal points in more

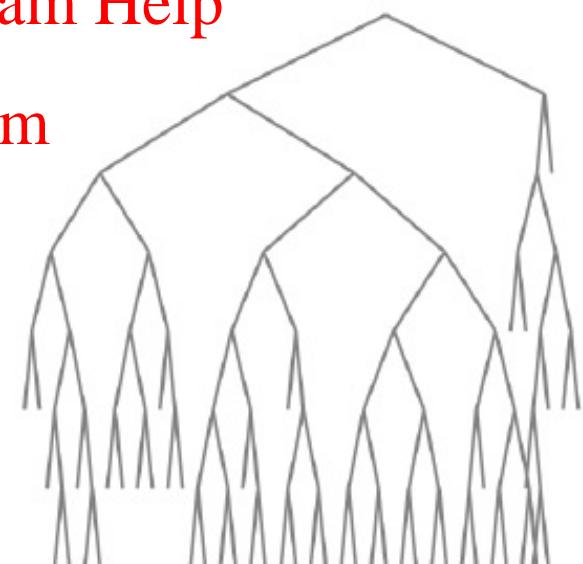


WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

- For each tree:
- Get a sample of the data
 - Randomly select a dimension
 - Randomly pick a value
- Draw a straight line through the data at that value and split data
- Repeat until tree is complete
- Generate multiple trees → forest
- Anomalies will be isolated in only a few steps
- Nominal points in more



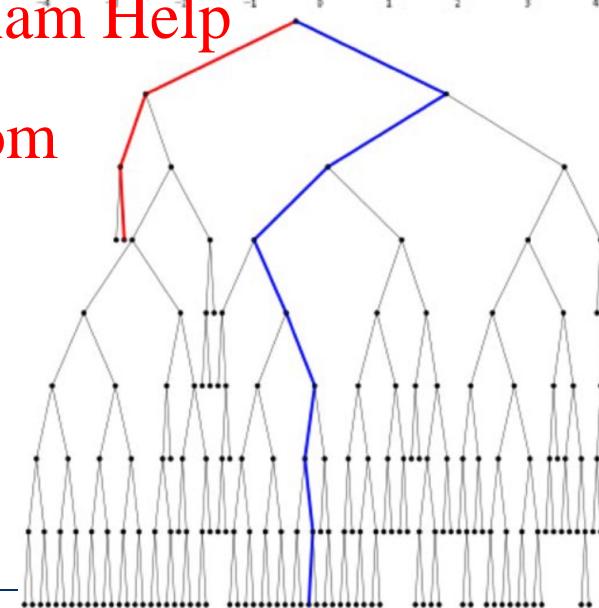
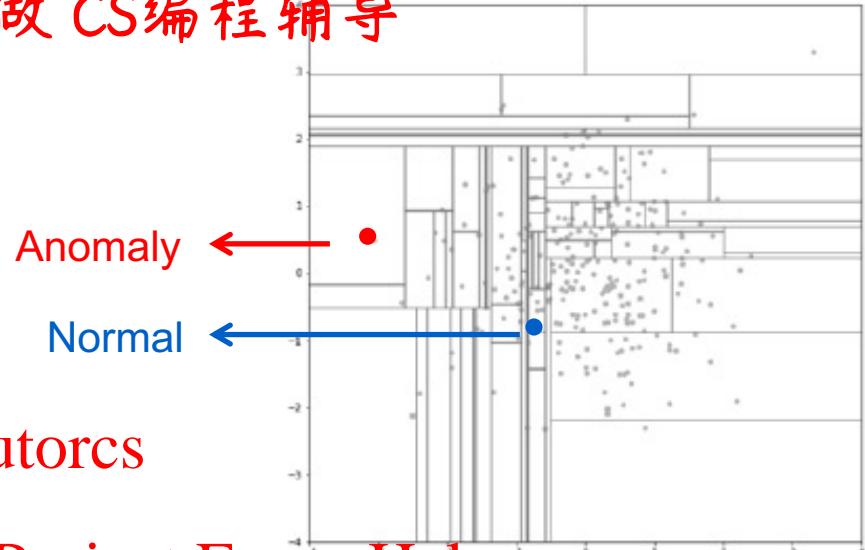
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



- Isolation Forest score:

程序代写代做 CS 编程辅导

$$\text{QR code} = 2^{-\frac{E(h(x))}{c(n)}}$$



- Where,

- $h(x)$ is the *path length* of observation x from the root node,
- $E(h(x))$ is the average of $h(x)$ from a collection of isolation trees
- n is the number of data points
- $c(n) = 2H(n - 1) - \left(\frac{2^{(n-1)}}{\pi}\right)$, where Euler's constant

Assignment Project Exam Help

Email: tutorcs@163.com

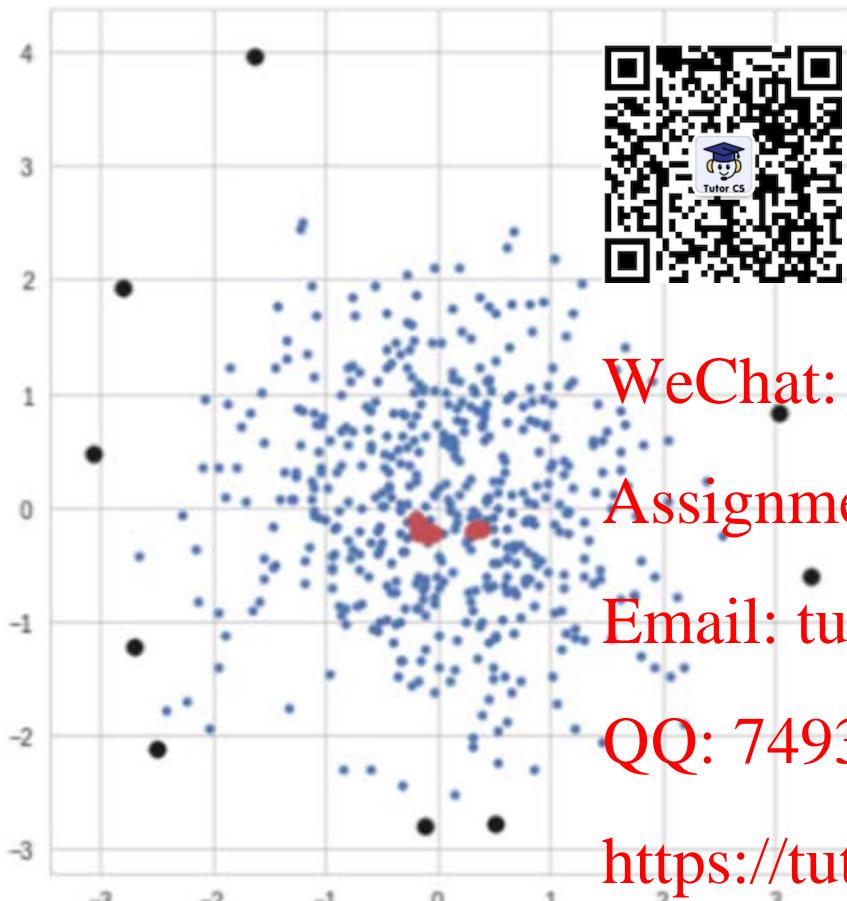
- $0 < s \leq 1$

QQ: 749389476

- $s \rightarrow 1$, then samples are definitely anomalies,
<https://tutorcs.com>
- $s \ll 0.5$, then samples are quite safe to be regarded as normal,
- $s = 0.5$, then the entire sample does not really have any distinct anomaly.

iForest Score – Case Study

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>



Advantages of iForest

程序代写代做 CS编程辅导

- Requires two parameters, the number of trees to build and the sub-sampling size
- Converges quickly with a small number of trees, and it only requires a small sub-sampling size to achieve high detection performance with high efficiency

WeChat: cstutorcs

- The isolation characteristic of iTrees enables them to *build partial models* and exploit sub-sampling to an extent that is not feasible in existing methods.

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

- Scales up to handle extremely large and high-dimensional datasets



程序代写代做 CS编程辅导

- What is anomaly detection and what are different types of anomalies?
- How we can evaluation the performance of anomaly detection techniques?
- How anomaly detection is different from other machine learning problems?
- How does the iForest algorithm operates, and what are its advantages of this method?



WeChat: cstutorcs

Assignment Project Exam Help

Next: Clustering and Density Based Anomaly Detection

Email: tutors@163.com

QQ: 749389476

<https://tutorcs.com>

- 程序代写代做 CS 编程辅导
1. Data Mining and Machine Learning in Security, Chapters 1,3.
 2. Machine Learning and Security, Chapter 1.
 3. Fei Tony Liu, Kai Ming Wong, Jiahu Zhou, “Isolation Forest”, IEEE International Conference on Data Mining, 2008.



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>