



FIT5202 – Data Processing for Big Data

Revision
Revised by
Chee-Ming Ting

Developed by
Prajwol Sangat



程序代写代做 CS 编程辅导



MONASH
INFORMATION
TECHNOLOGY

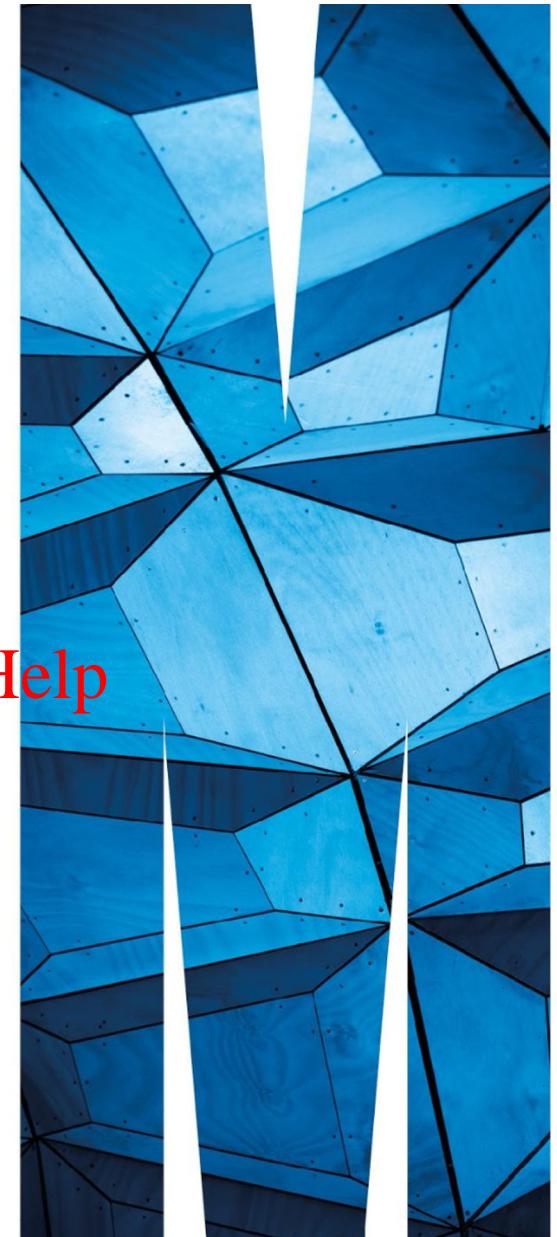
WeChat: esutores

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>





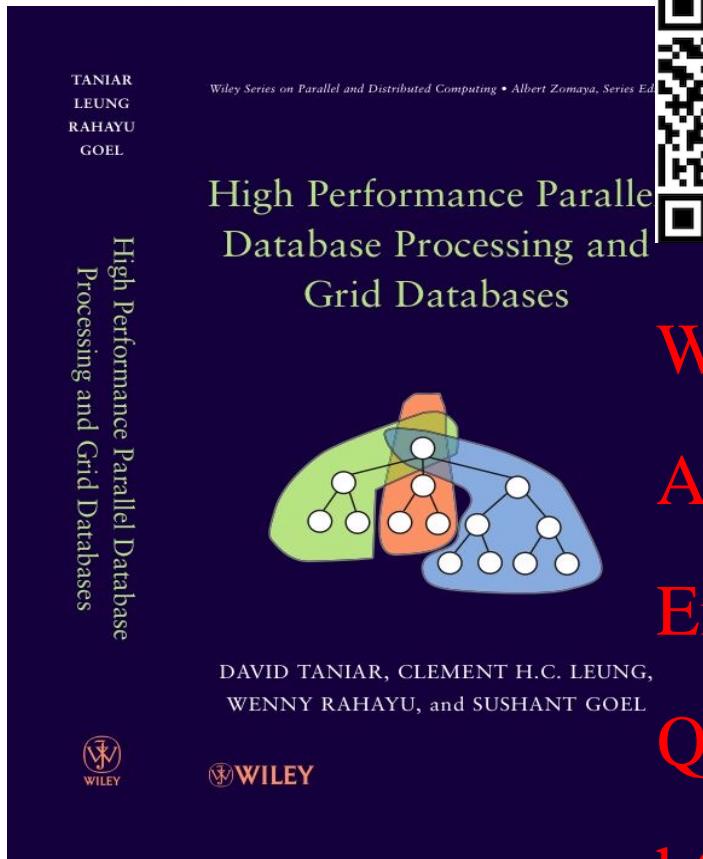
1. **Volume** → Sessions 1, 2, 3, 4
 - How to process Big Data of Large Volume?

2. **Complexity** → Sessions 5, 6, 7, 8
 - How to apply machine learning algorithms to every aspect of Big Data?

3. **Velocity** → Sessions 9, 10, 11
 - How to handle and process Fast Streaming Data?

<https://tutorcs.com>

Volume → Session 程序代写代做 CS编程辅导



Chapter 1 Introduction

1.1 A Brief Overview - Parallel Databases and Grid
WeChat: cstutorcs

1.2 Parallel Query Processing: Motivations

1.3 Parallel Query Processing: Objectives

1.4 Forms of Parallelism

1.5 Parallel Database Architectures

1.6 Grid Database Architecture

1.7 Structure of this Book

1.8 Summary

1.9 Bibliographical Notes

1.10 Exercises

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

1.3. Objectives (cont'd)



• Speed up

- Performance improvement gained by adding extra processing elements added
- Running a given task in less time by increasing the degree of parallelism

- Linear speed up: performance improvement growing linearly with additional resources
- Superlinear speed up
- Sublinear speed up

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

$$\text{Speed up} = \frac{\text{elapsed time on uniprocessor}}{\text{elapsed time on multiprocessors}}$$

<https://tutorcs.com>

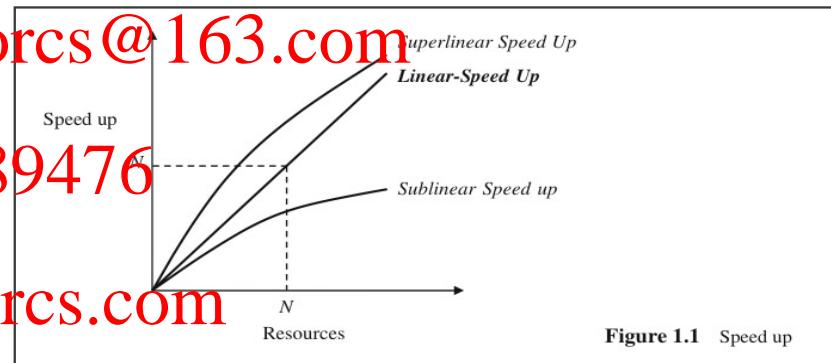


Figure 1.1 Speed up

程序代写代做 CS编程辅导



• Scale up

- Handling of larger tasks by increasing the degree of parallelism
 - The ability to process larger tasks in the same amount of time by providing more resources.
-
- Linear scale up: the ability to maintain the same level of performance when both the workload and the resources are proportionally added
 - Transactional scale up
 - Data scale up

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

$$\text{Scale up} = \frac{\text{uniprocessor elapsed time on small system}}{\text{multiprocessor elapsed time on larger system}}$$

<https://tutorcs.com>

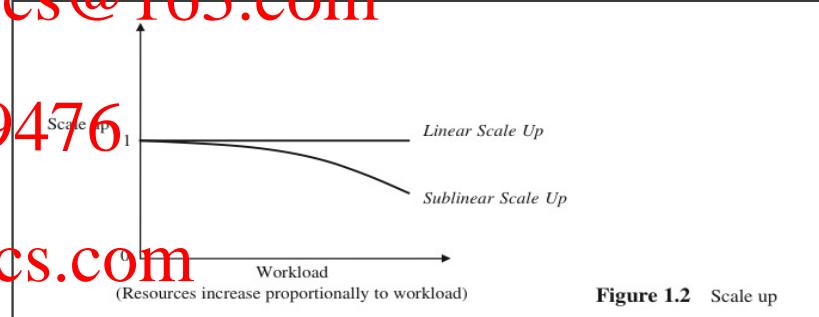


Figure 1.2 Scale up

Flux Quiz 1

程序代写代做 CS编程辅导

Using the current processing resources, we can finish processing 1TB (one terabyte) of data in 1 hour. Recently the amount of data has increased to 2TB and the management has decided to double the processing resources. Using the new processing resources, we can finish processing the 2 TB in 60 minutes.

Is this **speed up** or **scale up**? WeChat: cstutorcs

Solution:

Assignment Project Exam Help

$$\text{Scale up} = \frac{\text{uniprocessor elapsed time on small system}}{\text{mult processor elapsed time on larger system}}$$

Using x resources (current resources), 1TB = 60 minutes

When the resources are doubled (e.g. 1K becomes 2K now), a linear scale up is being able to complete 2TB in 60 minutes.

QQ: 749389476

In the question, using 2x resources, it finishes 2 TB in 60 minutes.

Therefore, it is linear scale up.

<https://tutorcs.com>

Parallel Obstacles

· Start-up and Consolidation

- Start up: initiation of multiple processes
- Consolidation: the cost for collecting results obtained from each processor by a host processor



WeChat: cstutorcs

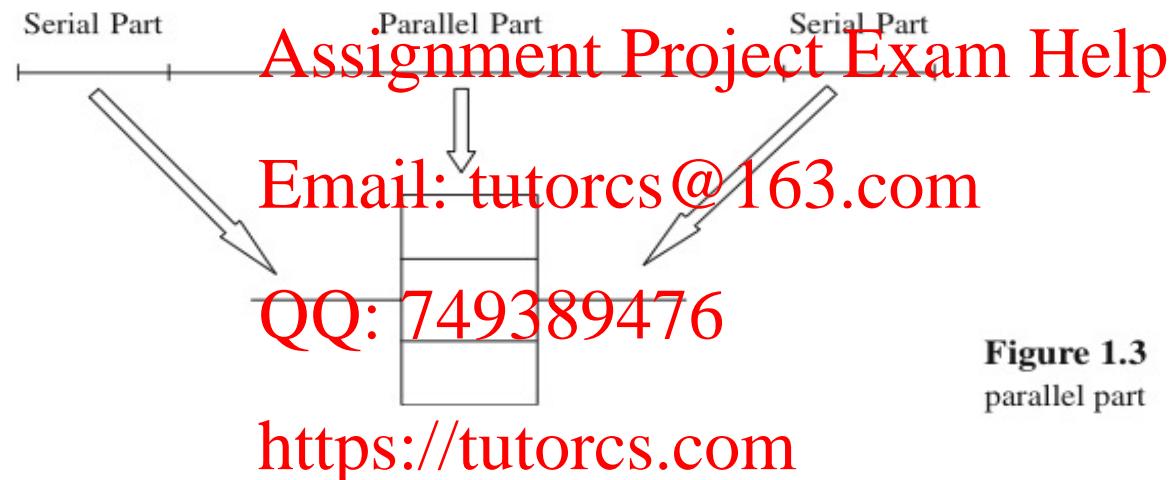


Figure 1.3 Serial part vs. parallel part

Parallel Obstacles

- #### **• Interference and Communication**

- Interference: competing to access resources
 - Communication: one process communicating with other processes, and often one has to wait for others to be ready for communication (i.e. waiting time).



WeChat: cstutorcs

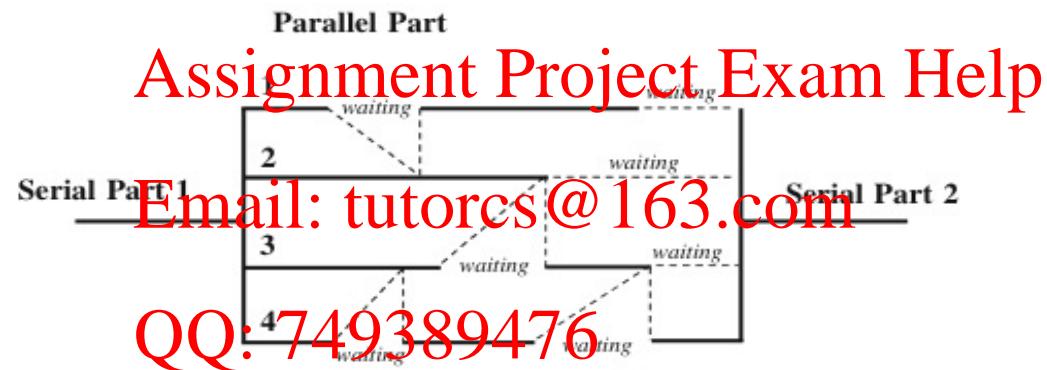


Figure 1.4 Waiting period
<https://tutorcs.com>

Flux Quiz 2

程序代写代做 CS编程辅导

There is a job that will take 60min to complete, if this is done by 1 processor. The serial part of this job is 10% of the total time. If we have 4 processors to use in this job, but each processor will have an overhead of 20% due to waiting time, communication time, etc. What type of speed up do we get?



Solution:

1 processor = 60min; Serial part = 10% = 6min; Parallel part = 54min

4 processors = $54\text{min}/4 = 13.5\text{min}$

Overhead = 20%

Hence, parallel processing part = $13.5\text{min} + 20\%\text{overhead} =$

$13.5\text{min} + 2.7\text{min} = 16.2\text{min}$

Total time = $6\text{min} + 16.2\text{min} = 22.2\text{min}$

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Speed up = $60\text{min} / 22.2\text{min} = 2.7$

Linear speedup should be 4. Speedup of 2.7 is Sublinear Speedup

$$\text{Speed up} = \frac{\text{elapsed time on uniprocessor}}{\text{elapsed time on multiprocessors}}$$

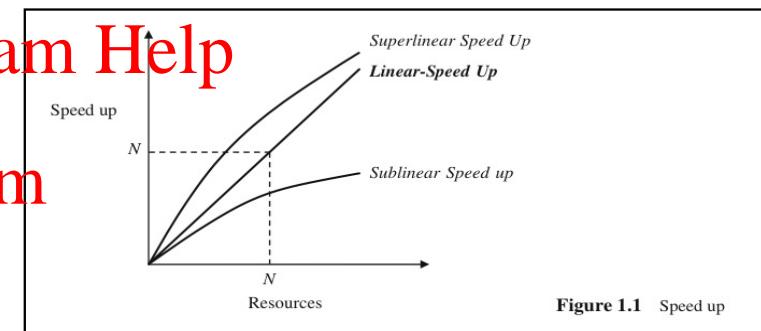


Figure 1.1 Speed up

程序代写代做 CS编程辅导

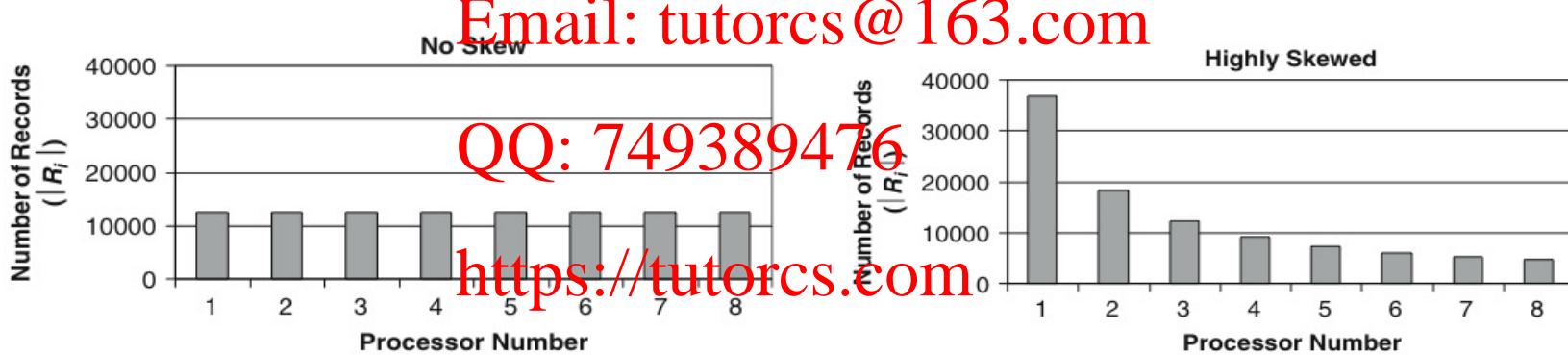
• Skew

- Zipf distribution model. Measured in terms of different sizes of fragments allocated to the processor.



$$= \frac{|R|}{i^\theta \times \sum_{j=1}^N \frac{1}{j^\theta}} \quad \text{where } 0 \leq \theta \leq 1 \quad (2.1)$$

- The symbol θ denotes the degree of skewness, where $\theta = 0$ indicates no skew, and $\theta = 1$ indicates highly skewed
- $|R|$ is number of records in the table, $|R_i|$ is number of records in processor i , and N is number of processor (i is a loop counter starting from 1 to N)
- Example: $|R|=100,000$ records, $N=8$ processors



Flux Quiz 3

程序代写代做 CS编程辅导

There 100,000 records in the table. These records are distributed to 32 processors. Assuming that the **skewness** degree is high, what is the estimated number of records in the heaviest processor?



$$|R_i| = \frac{|R|}{t^\theta \times \sum_{j=1}^N j^\theta} \quad \text{where } 0 \leq \theta \leq 1 \quad (2.1)$$

WeChat: cstutorcs

Assignment Project Exam Help

Solution:

i = 1 (heaviest processor)

$\theta = 1$

$$1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots + \frac{1}{32} = 4.05$$

$$\text{Number of records} = 100,000 / 4.05 = 24600$$

QQ: 749389476

<https://tutorcs.com>

Skew

程序代写代做 CS编程辅导

- Data Skew

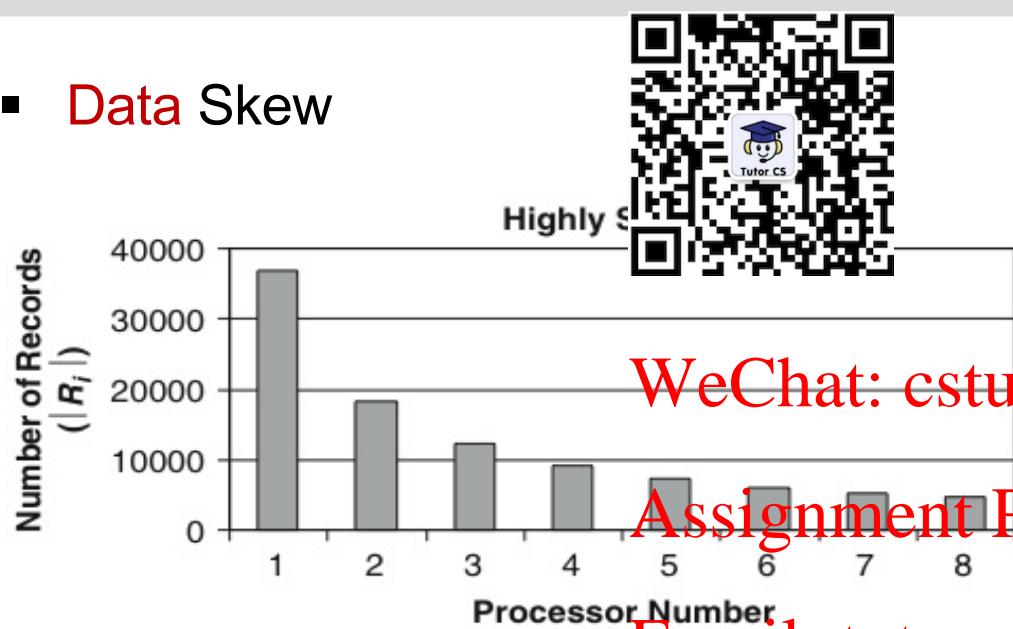
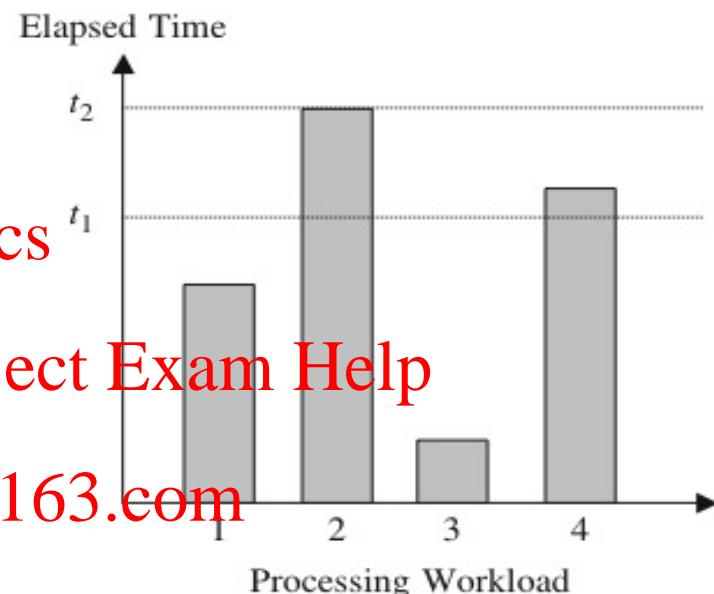


Figure 2.2 Highly skewed distribution

- Processing Skew



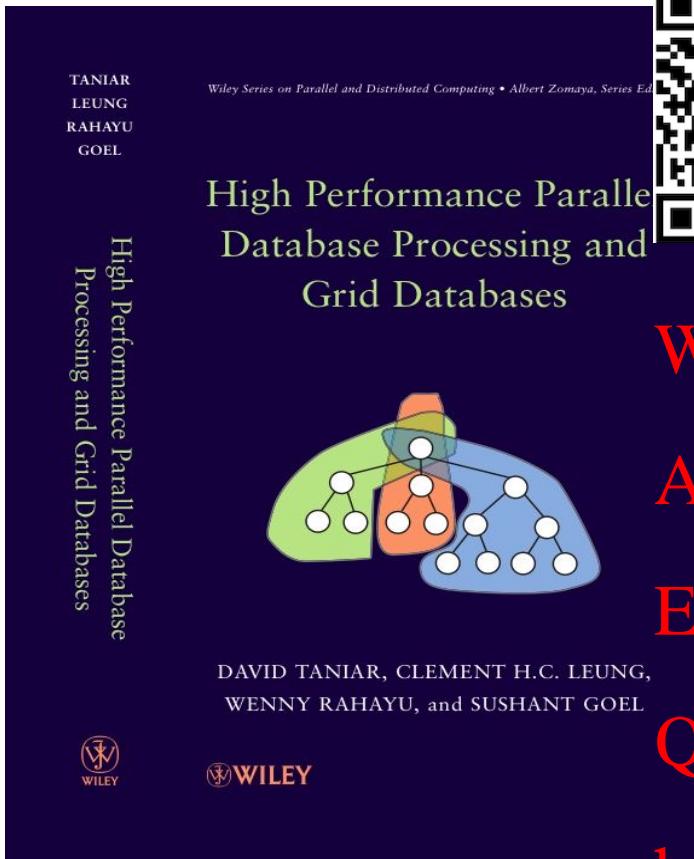
Uneven distribution of data in terms of size.

<https://tutorcs.com>

QQ: 749389476

Uneven processing time of the processors due to data skew.

Volume → Session 程序代写代做 CS编程辅导



Chapter 3 Parallel Search

WeChat: cstutorcs
Two steps
Data partitioning
- Parallel search

Assignment Project Exam Help

3.1 Search Queries

3.2 Data Partitioning

3.3 Search Algorithms

3.4 Summary

3.5 Bibliographical Notes

3.6 Exercises

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Basic Data Partitioning

- Round-robin or random data partitioning
- Hash data partitioning
- Range data partitioning
- Random-unequal data partitioning

WeChat: cstutorcs

Complex Data Partitioning

Assignment Project Exam Help

- Complex data partitioning is based on multiple attributes or is based on a single attribute but with multiple partitioning methods
- Hybrid-Range Partitioning Strategy (HRPS)
- Multiattribute Grid Declustering (MAGIC)
- Bubba's Extended Range Declustering (BERD)

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Data Partitioning (cont'd)

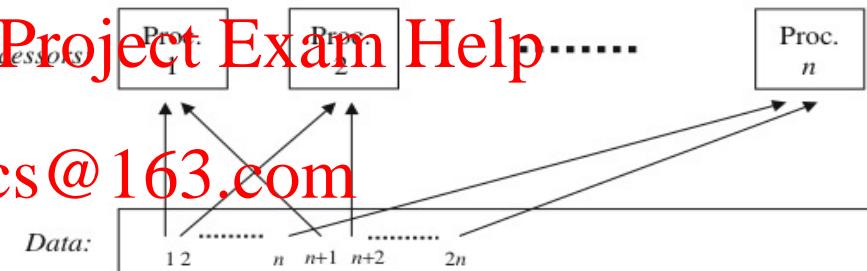


Round-robin data partitioning

- Each record in turn assigned to a processing element in a clockwise manner
- “Equal partitioning” or “Random-equal partitioning”
- Data evenly distributed, hence supports load balance
- But data is not grouped semantically

WeChat: csutorcs

Assignment Project Exam Help.....



Email: tutorcs@163.com

QQ: 749389476

Figure 3.3 Round-robin data partitioning

<https://tutorcs.com>

Data Partitioning (cont'd)

Hash data partitioning

- A hash function is used to partition the data
- Hence, data is grouped logically, that is data on the same group shared the same hash value
- Selected processors may be identified when processing a search operation (exact-match search), but for range search (especially continuous range), all processors must be used
- Initial data allocation is not balanced either

WeChat: csutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

The data is divided using a hash function
can cause skewness

<https://tutorcs.com>

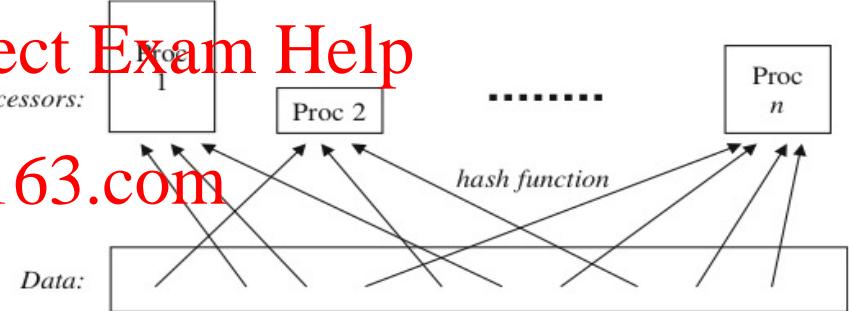


Figure 3.4 Hash data partitioning

程序代写代做 CS编程辅导

Data Partitioning (cont')

Range data partitioning

- Spreads the record across all processors
- Processing records in a given range of the partitioning attribute
- A specific range can be directed to certain processor only
- Initial data allocation is skewed too



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

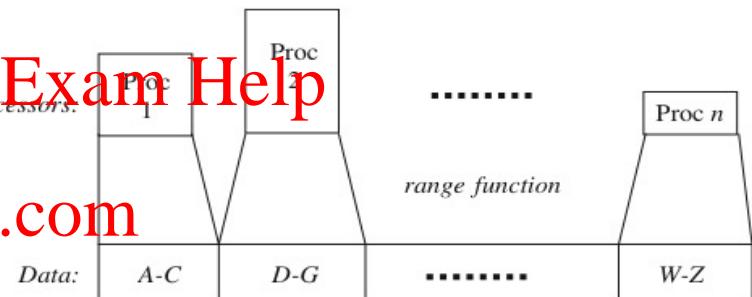


Figure 3.5 Range data partitioning

Search Algorithms

程序代写代做 CS编程辅导

Serial search algorithm

- Linear search
- Binary search



Parallel search algorithms:

- Processor activation or involvement
- Local searching method
- Key comparison

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Search Algorithms (cont.)

Processor activation or involvement

- The number of processors involved in the search depends on the data partitioning method used by the algorithm
- If we know where the data items sought are stored, then there is no point in activating all other processors in the searching process
- Depends on the data partitioning method used
- Also depends on what type of selection query is performed



Table 3.6 Processor activation or involvement of parallel search algorithms

		Data Partitioning Methods			
		Random-Equal	Hash	Range	Random-Inequal
Exact Match		All	1	1	All
Range Selection	Continuous	All	All	Selected	All
Selection	Discrete	All	Selected	Selected	All

Email: tutorcs@163.com

QQ: 749389476

How many processors we need to activate to do the search?

<https://tutorcs.com>

Search Algorithms (cont.)

Local searching methods

- The searching range is limited to the processor(s) involved in the searching process
- Depends on the data ordering, regarding the type of the search (exact match or range)



WeChat: cstutorcs

Table 3.7 Local searching method of parallel search algorithms

		Records Ordering	
		Ordered	Unordered
Exact Match	Continuous	Binary Search	Linear Search
Range	Continuous	Binary Search	Linear Search
Selection	Discrete	Binary Search	Linear Search

QQ: 749389476

Key idea: if the data is sorted, we use binary.

If not we used linear.

<https://tutorcs.com>

程序代写代做 CS编程辅导

Search Algorithms (con)

Key comparison

- Compares the data in the table with the condition specified by the query
- When a match is found: continue to find other matches, or terminate
- Depends on whether the data in the table is unique or not



Assignment Project Exam Help

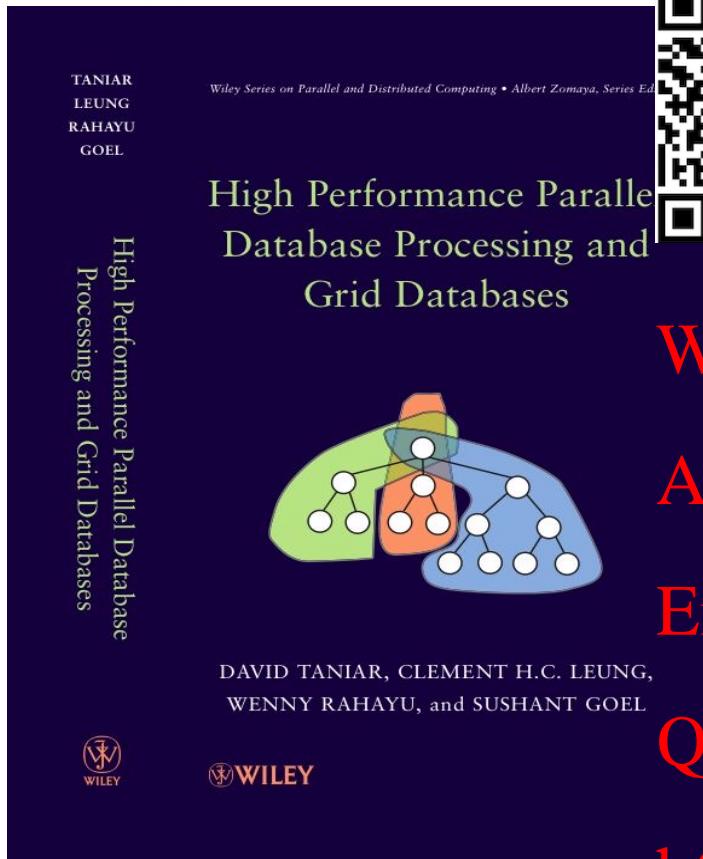
		Search Attribute Values	
		Unique	Duplicate
Exact Match	Stop	Stop	Continue
	Continuous	Continue	Continue
	Discrete	Continue	Continue

QQ: 749389476

Should we stop the searching when the match is found?

<https://tutorcs.com>

Volume → Session 程序代写代做 CS编程辅导



Chapter 5 Parallel Join

- WeChat: cstutorcs
Assignment Project Exam Help
Email: tutorcs@163.com
QQ: 749389476
- 5.1 Join Operations
5.2 Serial Join Algorithms
5.3 Parallel Join Algorithms
5.4 Cost Models
5.5 Parallel Join Optimization
5.6 Summary
5.7 Bibliographical Notes
5.8 Exercises

<https://tutorcs.com>

Join Algorithms

程序代写代做 CS编程辅导



- Parallel Inner Join com
 - **Data Partitioning**
 - Divide and Broadcast
 - Disjoint Partitioning
 - **Local Join**
 - Nested-Loop Join
 - Sort-Merge Join
 - Hash Join
- Example of a Parallel Inner Join Algorithm
 - **Divide and Broadcast, plus Hash Join**

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Serial Join Algorithms



Hash-based Join Algorithms

- The records of files R and S are both hashed to the *same hash file*, using the *same hashing function* on the join attributes A of R and B of S as hash keys
- A single pass through the file with fewer records (say, R) hashes its records to the hash file buckets
- A single pass through the other file (S) then hashes each of its records to the appropriate bucket, where the record is combined with all matching records from R

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

<https://tutorcs.com>

程序代写代做 CS编程辅导

Serial Join Algorithms



Table S	
Arts	8
Business	15
CompSc	2
Dance	11
Engineering	7
Finance	21
Geology	10
Health	11
IT	18

Hash Table	
Index	Entries
1	Geology/10
2	CompSc/2
3	Dance/12
4	
5	
6	Business/15
7	Engineering/7
8	Arts/8
9	IT/18
10	
11	
12	

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Figure 5.6. Hashing Table S

<https://tutorcs.com>

程序代写代做 CS编程辅导

Serial Join Algorithms (c)

Table R	
Adele	8
Bob	22
Clement	16
Dave	23
Ed	11
Fung	25
Goel	3
Harry	17
Irene	14
Joanna	2
Kelly	6
Lim	20
Meng	1
Noor	5
Omar	19



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Figure 5.7. Probing Table R

<https://tutorcs.com>

程序代写代做 CS编程辅导

Parallel Join Algorithm



Divide and Broadcast-based Parallel Join Algorithms

- Two stages: data partitioning using the divide and broadcast method, and a local join
- Divide and Broadcast method: Divide one table into multiple disjoint partitions, where each partition is allocated a processor, and broadcast the other table to all available processors
- Dividing one table can simply use equal division
- Broadcast means replicate the table to all processors
- Hence, choose the smaller table to broadcast and the larger table to divide

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

Parallel Join Algorithms



WeChat: cstutorcs

Assignment Project Exam Help

Processor S1		Processor S2		Processor S3	
R1		R2		R3	
Adele	8	Arts	8	Business	12
Bob	22	Dance	15	Engineering	7
Clement	16	Fung	25	Health	11
Dave	23	Goel	3		
Ed	11	Harry	17		
		Jess	14		
		Wanida	12		
				Kelly	6
				Lim	20
				Meng	1
				Noor	5
				Omar	19
				CompSc	2
				Finance	21
				IT	18

Email: tutorcs@163.com

QQ: 749389476

Figure 5.10 Initial data placement

<https://tutorcs.com>

程序代写代做 CS编程辅导

Parallel Join Algorithms



Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Figure 5.11 Divide and broadcast result
<https://tutorcs.com>

程序代写代做 CS编程辅导

Parallel Join Algorithms



Processor 1

Result 1

Adele	8	Arts
Ed	11	Health

Processor 2

Result 2

Joanna	2	CompSc
--------	---	--------

Processor 3

Result 3

NIL

WeChat: cstutorcs

Processor 1

R1

Adele	8	S1
Bob	22	Arts
Clement	16	Dance
Dave	23	Geology
Ed	11	Business

S2

Business	12
Engineering	7
Health	11

S3

CompSc	2
Finance	21
IT	18

Processor 2

R2

Ed	25	S2
Karen	3	Business
Harry	17	Engineering
Irene	14	Health
Joanna	2	IT

S1

Arts	8
Dance	15
Geology	10

S3

CompSc	2
Finance	21
IT	18

Processor 3

R3

Kelly	6	S3
Lily	10	CompSc
Meng	1	Finance
Noor	5	IT
Omar	19	IT

S1

Arts	8
Dance	15
Geology	10

S2

Business	12
Engineering	7
Health	11

Assignment Project Exam Help

Email: tutorcs@163.com
QQ: 749389476

Figure 5.12 Join results based on divide and broadcast

<https://tutorcs.com>

Cost Models?

程序代写代做 CS编程辅导



Divide and Broadcast

- Join two tables **R** and **table S**)
- The two tables have been partitioned and stored in 3 processors **WeChat: cstutorcs**
- The tables have been partitioned using the random-equal data partitioning method **Assignment Project Exam Help**
- The table fragments are called R1, R2, R3, and S1, S2, S3 (in general, each fragment is called **R_i** or **S_i**, where i is the processor number) **Email: tutorcs@163.com** **QQ: 749389476**

<https://tutorcs.com>

Flux Quiz 4

程序代写代做 CS编程辅导

$|S| = 600$ records, each record has a length of 100 bytes, and $N=3$.

Solution:

$S = 60000$ bytes

$Si = S/N = 60000/3 = 20000$ bytes

$|S| = 600$ records

$|Si| = 600/3 = 200$ records



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Transfer cost = $(Si/P) \times (N-1) \times (mp+ml)$

Receiving cost = $(S/P - Si/P) \times (mp)$

<https://tutorcs.com>

Table 2.1 Cost notations

Symbol	Description
Data parameters	
R	Size of table in bytes
R_i	Size of table fragment in bytes on processor i
$ R $	Number of records in table R
$ R_i $	Number of records in table R on processor i
Systems parameters	
N	Number of processors
P	Page size
H	Hash table size
Query parameters	
π	Projectivity ratio
σ	Selectivity ratio
Time unit cost	
IO	Effective time to read a page from disk
t_r	Time to read a record in the main memory
t_w	Time to write a record to the main memory
t_d	Time to compute destination
Communication cost	
mp	Message protocol cost per page
ml	Message latency for one page

Cost Models for Parallel Join



Cost Models for Divide and Broadcast

- Phase 1: data loading
 - . Scan cost for loading data in each processor is: $(S_i / P) \times IO$
 - . Select cost for getting record page is: (read and written by CPU to memory): $|S_i| \times (tr + tw)$

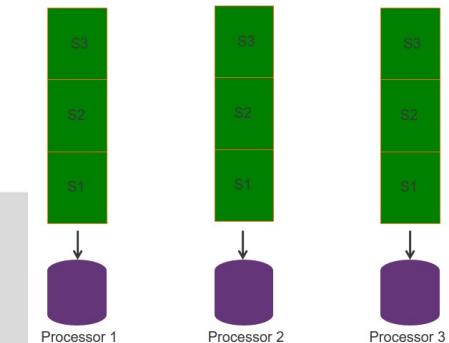
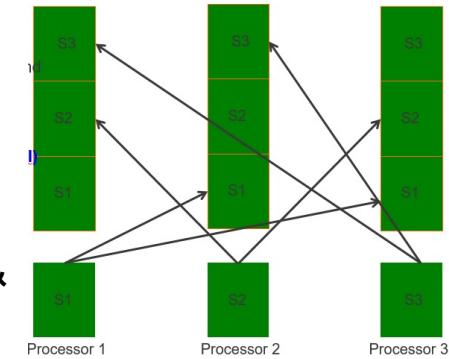
WeChat: cstutorcs

Assignment Project Exam Help
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

33



Query Parame

- **Projectivity ratio**
 - Ratio between attribute size and original record length
- **Selectivity ratio**
 - Ratio between number of records in the query result and original total number of records



WeChat: cstutorcs

Example: Join selectivity ratio

Assignment Project Exam Help

If the query operation involves two tables (like in a join operation), a selectivity ratio can be written as σ_j . For example, the value of σ_j indicates the ratio between the number of records produced by a join operation and the number of records of the Cartesian product of the two tables to be joined. For example, $|R_i| = 1000$ records and $|S_i| = 500$ records; if the join produces 5 records only, then the join selectivity ratio σ_i is $5/(1,000 \times 500) = 0.00001$.

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

Parallel Join Query



Parallel Outer Join Processing methods

ROJA (Redistribution Outer Join Algorithm)

DOJA (Duplication Outer Join Algorithm)

DER (Duplication & Efficient Redistribution)

Assignment Project Exam Help

Load Balancing

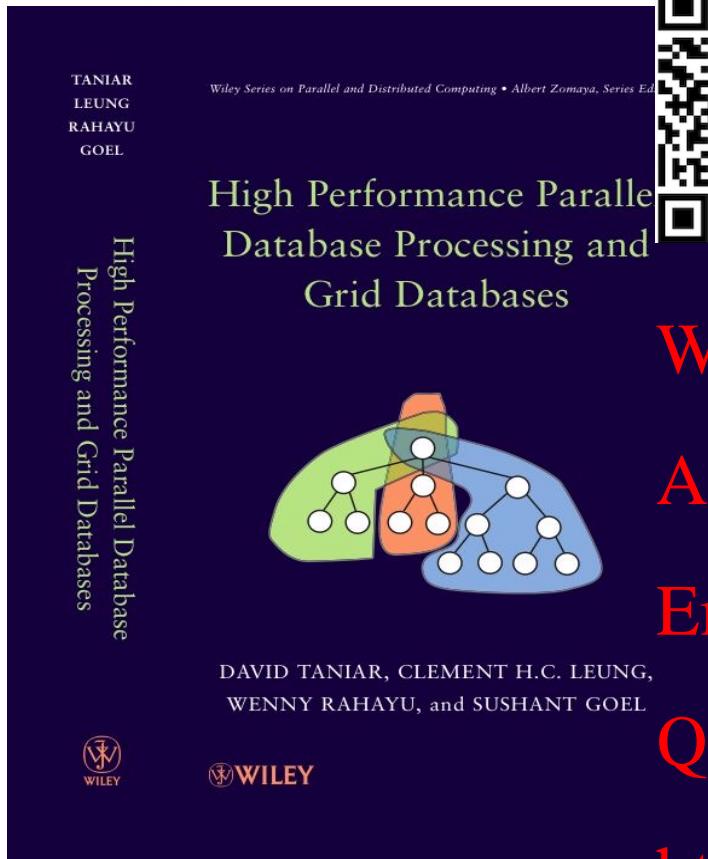
OJSO (Outer Join Skew Optimization)

QQ: 749389476

<https://tutorcs.com>

	ROJA	DOJA	DER
Steps	<p>Step 1: Distribute or reshuffle the data based on the join attribute.</p> <p>Step 2: Each processor performs the Local outer Join.</p>	<p>Step 1: Replication. We duplicate the small table.</p> <p>Step 2: Local Inner Join</p> <p>Step 3: Distribute the inner join results to each processor X.</p> <p>Step 4: Local outer join</p>  <p>WeChat: cstutorcs</p>	<p>Step 1: Replication. We broadcast the left table.</p> <p>Step 2: Local Inner Join</p> <p>Step 3: Select the ROW ID of left table with no matches.</p> <p>Step 4: Redistribute the ROW ID.</p> <p>Step 5: Store the ROW ID that appears as many times as the number of processors.</p> <p>Step 6: Inner join</p>
Pros	fast performance, only two steps	None. ROJA is faster than DOJA.	Redistributes dangling row IDs instead of actual records.
Cons	redistribution of data -> data skew, communication cost	<p>In the replication step, if the table is large, the replication cost is expensive.</p> <p>In the distribution step, data skew and communication cost similar to ROJA</p> <p>QQ: 749389476</p> <p>https://tutorcs.com</p>	In the replication step, if the table is large, the replication cost is expensive.

Volume → Session 程序代写代做 CS编程辅导



WeChat: cstutorcs

Chapter 4

Parallel Sort and GroupBy

Assignment Project Exam Help

4.1 Sorting, Duplicate Removal and Aggregate

4.2 Serial External Sorting Method

4.3 Algorithms for Parallel External Sort

4.4 Parallel Algorithms for GroupBy Queries

4.5 Cost Models for Parallel Sort

4.6 Cost Models for Parallel GroupBy

4.7 Summary

4.8 Bibliographical Notes

QQ: 749389476

<https://tutorcs.com>

Sorting, and Serial Sorting



- Serial Sorting – **INTERNAL**
 - The data to be sorted FITS entirely into the main memory
 - Bubble Sort
 - Insertion Sort
 - Quick Sort
- Serial Sorting - **EXTERNAL**
 - The data to be sorted DOES NOT fit entirely into the main memory
 - Sort-Merge

WeChat: cstutorcs

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Flux Quiz 5

程序代写代做 CS编程辅导

There are 150 data pages to be sorted. The machine that we have has a limited memory, and can only take 8 pages at a time.

How many passes will it take to sort 150 data pages?



Solution:

File size to be sorted = 150 pages, number of buffer (or memory size) = 8 pages

Number of subfiles = $150/8 = 19$ subfiles (Last subfile has only 6 pages)

Pass 0 (sorting phase): For each subfile, read from disk, sort in main-memory, and write to disk

WeChat: csstutorcs
Assignment Project Exam Help

Merging phase: We use 7 buffers for input and 1 buffer for output

Pass 1: Read 7 sorted subfiles and perform 7-way merging. Repeat the 7-way merging until all subfiles are processed. Result = 3 subfiles

Pass 2: Merge the 3 subfiles

QQ: 749389476

Summary: 150 pages and 8 buffer pages require 3 passes

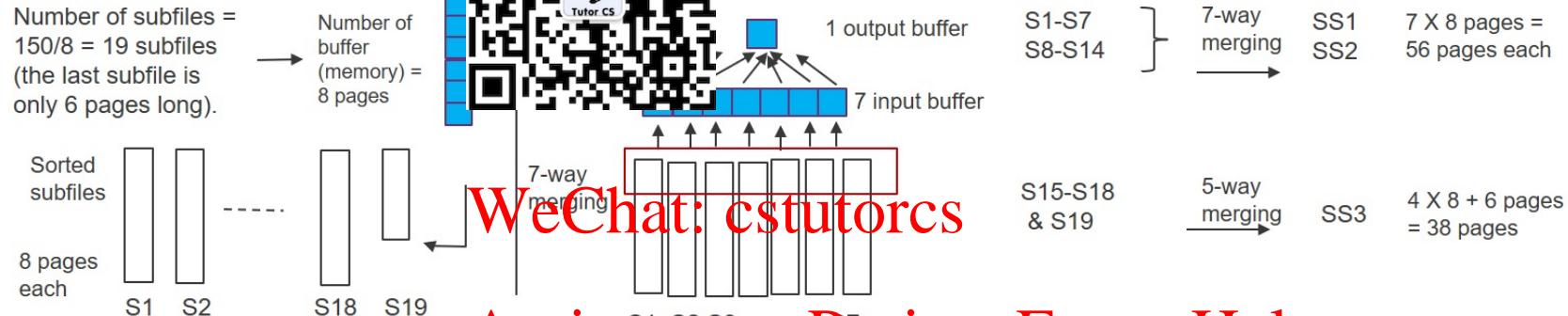
<https://tutorcs.com>

程序代写代做 CS编程辅导

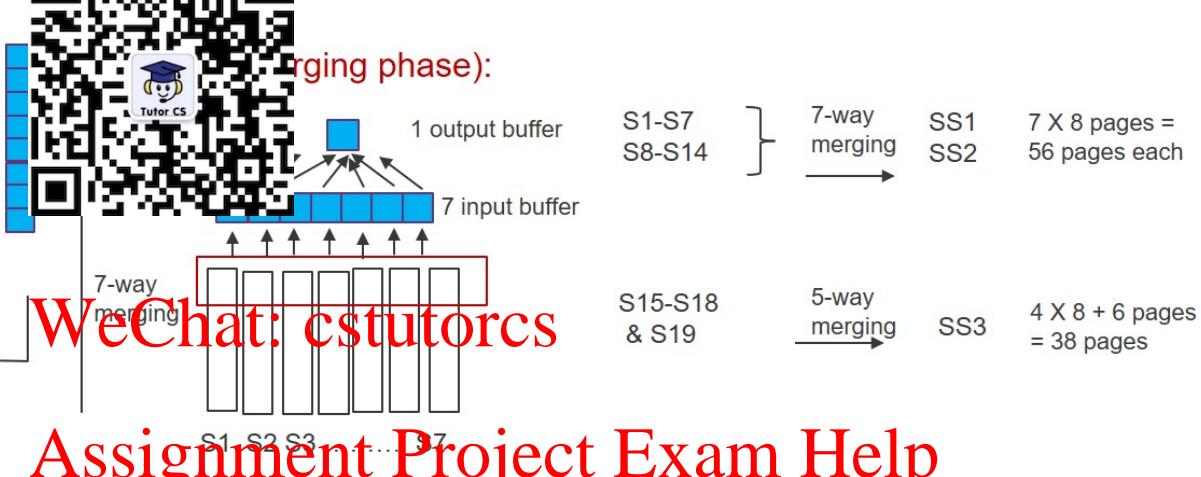
File size to be sorted = 150 pages of buffer (or memory size) = 8 pages

Pass 0 (sorting phase):

Number of subfiles =
 $150/8 = 19$ subfiles
(the last subfile is
only 6 pages long).



Merging phase:



WeChat: cstutorcs

Assignment Project Exam Help
Pass 2 (Merging phase):
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Parallel External Sort

程序代写代做 CS编程辅导

- Parallel Merge-All Sort
- Parallel Binary-Merge Sort
- Parallel Redistribution Binary-Merge Sort
- Parallel Redistribution Merge-All Sort
- Parallel Partitioned Sort



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Parallel External Sort (c)



Parallel Merge-All Sort

- A traditional approach
- Two phases: local sort and final merge
- Load balanced in local sort
- Problems with merging:
 - Heavy load on one processor
 - Network contention

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

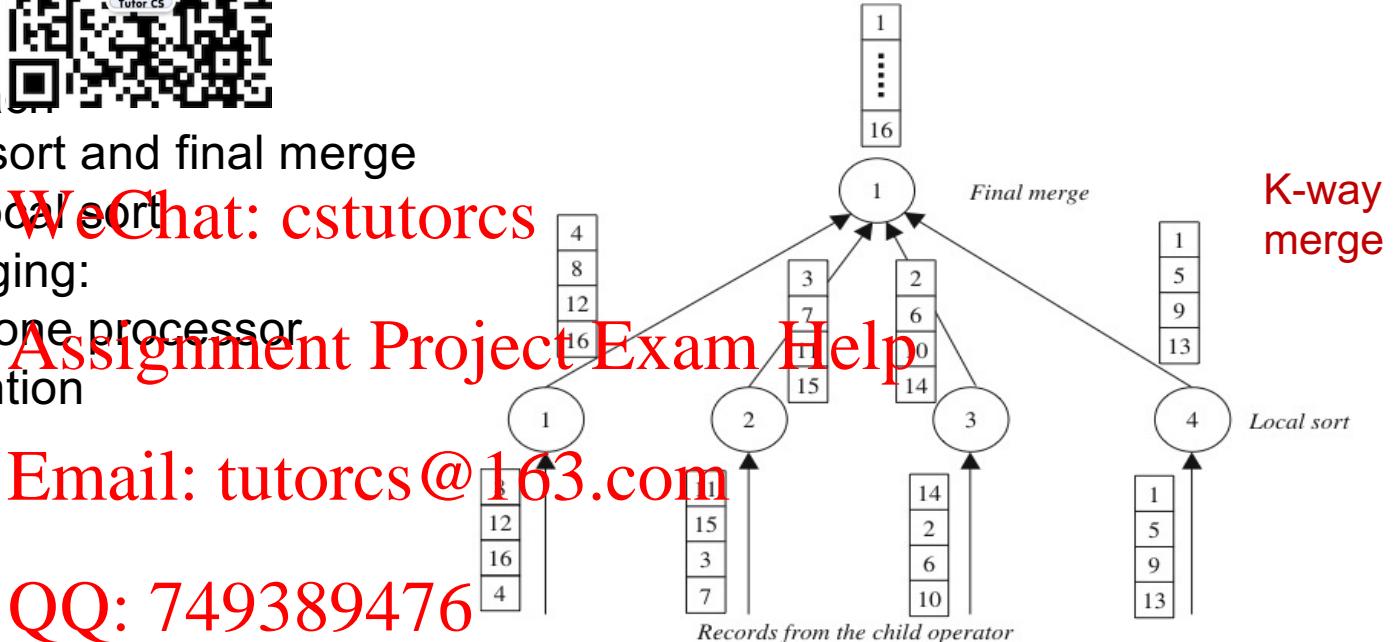
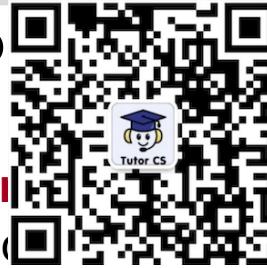


Figure 4.3 Parallel merge-all sort

Parallel External Sort (co

- Parallel Binary-Merge Sort
 - Local sort similar to traditional method
 - Merging in pairs only
 - Merging work is now spread to pipeline of processors, but merging is still heavy



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Figure 4.4 Parallel binary-merge sort

程序代写代做 CS编程辅导

Parallel Redistribution Binary-Merge Sort



Parallelism at all levels in the pipeline hierarchy

Step 1: local sort

Step 2: redistribute the results of local sort

Step 3: merge using the same pool of processors

Benefit: merging becomes lighter than without redistribution

Problem: height of the tree

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

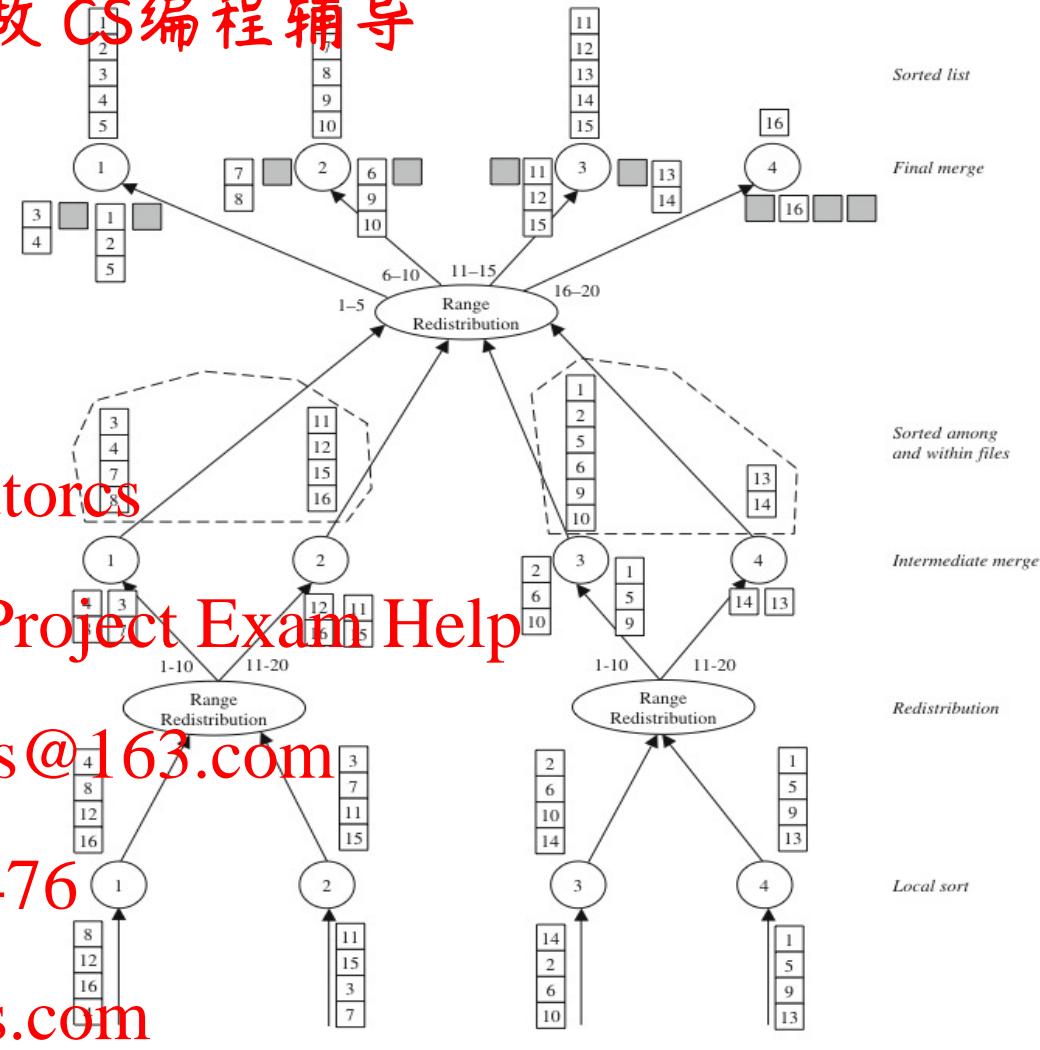


Figure 4.6 Parallel redistribution binary-merge sort

程序代写代做 CS编程辅导

Parallel Redistribution Merge Sort

- Reduce the height of the tree while still maintaining parallelism
- Like parallel merge-all sort, but with redistribution
- The advantage is true parallelism in merging
- Skew problem in the merging

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

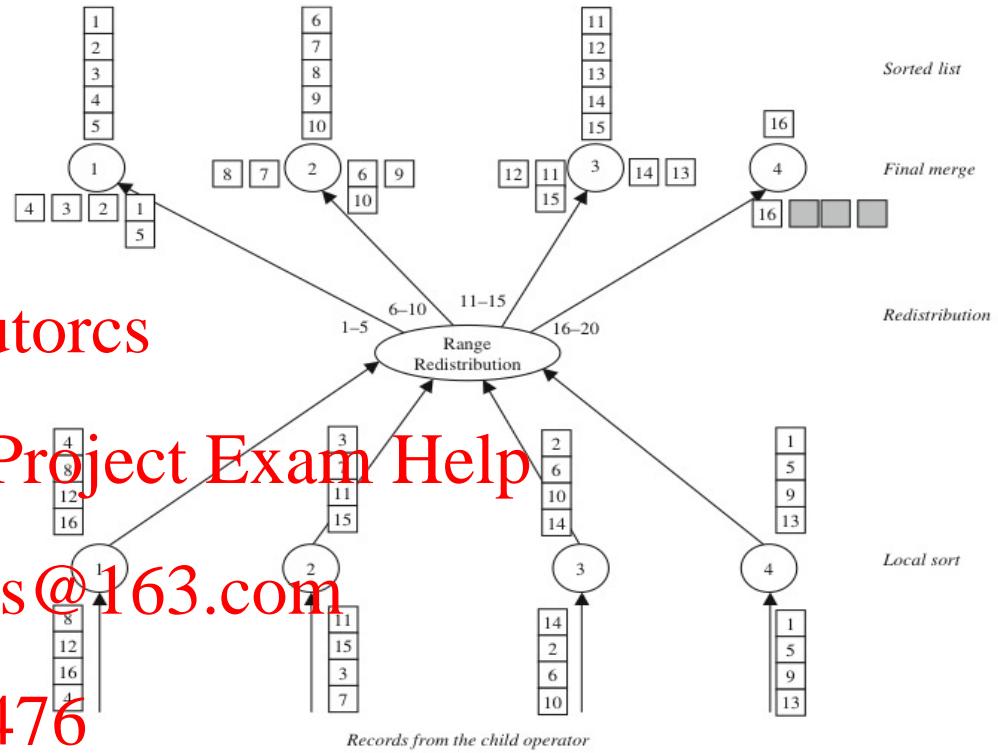


Figure 4.7 Parallel redistribution merge-all sort

程序代写代做 CS编程辅导

Parallel Partitioned Sort

- Two stages: Partitioning
Independent local work
- Partitioning (or range redistribution) may raise load skew
- Local sort is done after the partitioning, not before
- No merging is necessary
- Main problem: **Skew** produced by the partitioning



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

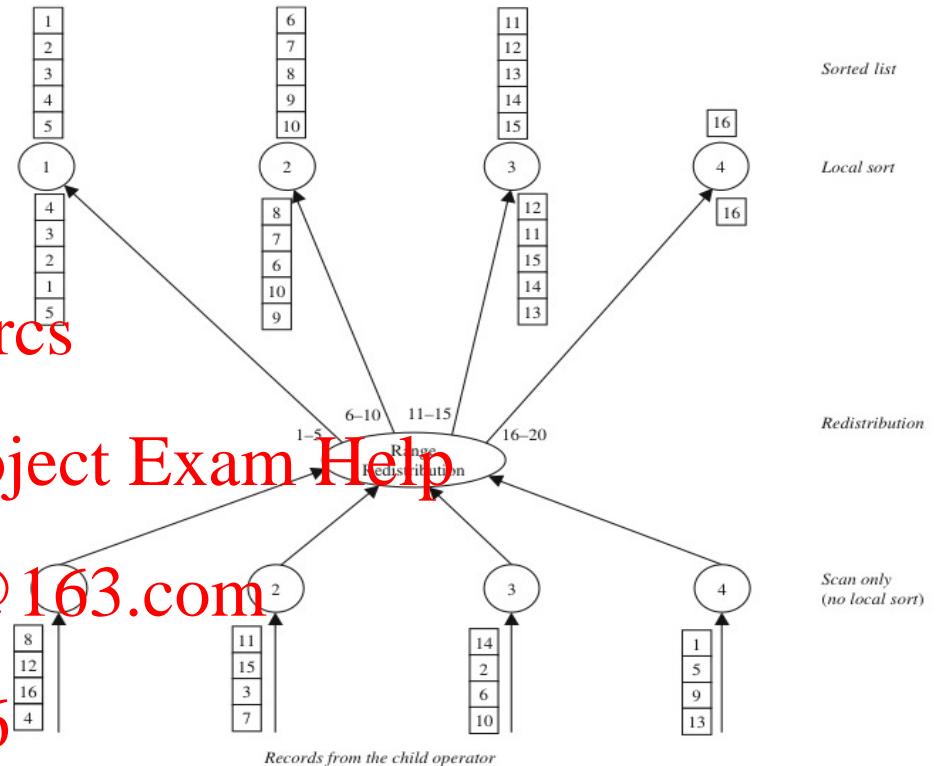


Figure 4.8 Parallel partitioned sort



Exercise

- Given a data set $D = \{8, 1, 12, 15, 2, 16, 3, 6, 9, 4, 7, 10, 13\}$ and four processors, show step by step how the **Parallel Partitioned Sort** works.

WeChat: cstutorcs

Initial Data Partitioning (Round Robin):

- $P_1 = \{8, 12, 16, 4\}$, $P_2 = \{11, 15, 3, 7\}$, $P_3 = \{14, 2, 6, 10\}$ and $P_4 = \{1, 5, 9, 13\}$

Range Redistribution (Range Logic: $P_1 = 1-5$, $P_2 = 6-10$, $P_3 = 11-15$, $P_4 = 16-20$)

- $P_1 = \{4, 3, 2, 1, 5\}$, $P_2 = \{8, 7, 6, 10, 9\}$, $P_3 = \{12, 11, 15, 14, 13\}$, $P_4 = \{16\}$

Local Sort

- $P_1 = \{1, 2, 3, 4, 5\}$, $P_2 = \{6, 7, 8, 9, 10\}$, $P_3 = \{11, 12, 13, 14, 15\}$, $P_4 = \{16\}$

<https://tutorcs.com>

Parallel Group By

程序代写代做 CS编程辅导

- Traditional methods (Map Reduce, Hash-based, Hierarchical Merging)
- Two-phase method
- Redistribution method



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Parallel Group By (cont'd)

Traditional Methods

Step 1: local aggregate processor

Step 2: global aggregation processor

May use a Merge-All or Hierarchical method

Need to pay a special attention to some aggregate functions (AVG) when performing a local aggregate process



WeChat: cstutors

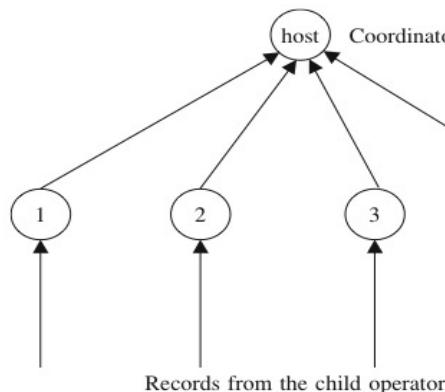


Figure 4.10 Traditional method

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

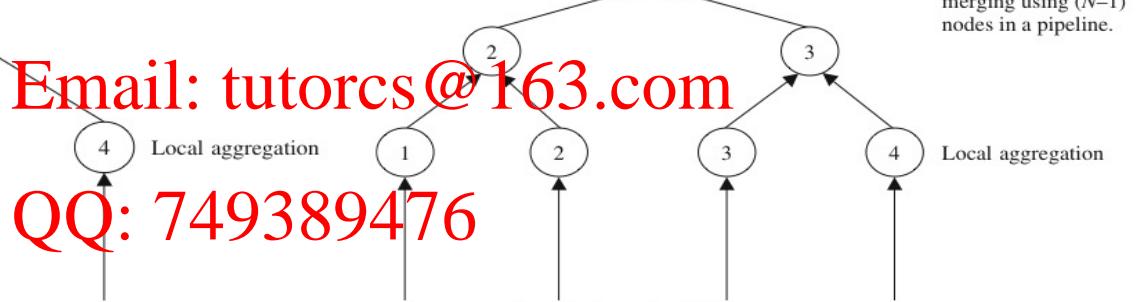


Figure 4.11 Hierarchical merging method

Two-level hierarchical merging using $(N-1)$ nodes in a pipeline.

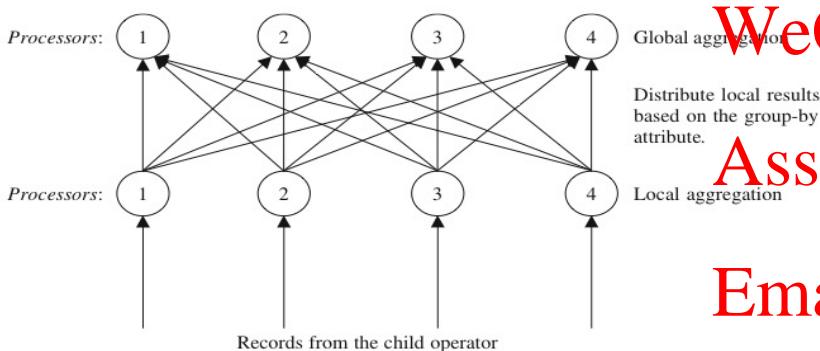
程序代写代做 CS编程辅导

Parallel Group By (cont'd)

Two-Phase Method

Step 1: local aggregate in each processor. Each processor groups local records according to the group-by attribute.

Step 2: global aggregation. Aggregation results from each processor are redistributed and then final aggregate is performed in each processor



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Final results			
Clayton 2	Hawthorn 1	Malvern 3	Richmond 2
Balwyn 1	Kew 0		Vermont 2
Elwood 2			
Caulfield 1			

Processor 1 Processor 2 Processor 3 Processor 4

Global Aggregation Phase

Distribute Local Aggregation Results Phase			
Clayton 2	Hawthorn 1	Malvern 1	Richmond 1
Balwyn 1	Kew 1	Malvern 1	Vermont 1
Elwood 1	Caulfield 1	Richmond 1	Vermont 1
Balwyn 1	Elwood 1	Malvern 1	Vermont 1
Clayton 1	Kew 1	Clayton 1	Vermont 1

Processor 1 Processor 2 Processor 3 Processor 4

Parallel Group By (cont'd)

Redistribution Method

Step 1 (Partitioning phase): distribute raw records to all processors

Step 2 (Aggregation phase): each processor performs a local aggregation

WeChat: cstutorcs

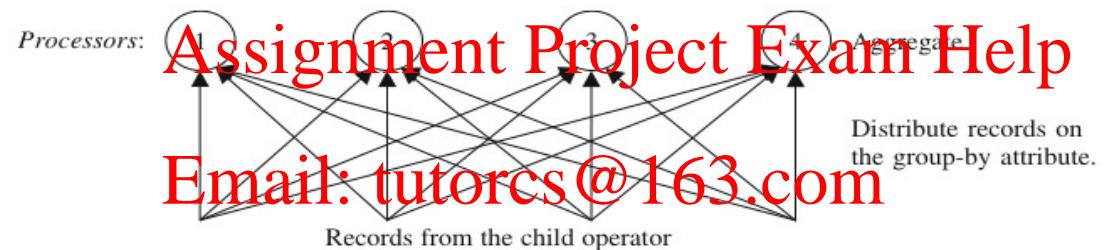


Figure 4.13 Redistribution method
QQ: 749389476

<https://tutorcs.com>

Flux Quiz 7

程序代写代做 CS编程辅导

The **Redistribution Method** has a task stealing option, through the Task Stealing method.



The **Two-Phase Method** does not have a task balancing problem.

Solution: False

WeChat: cstutorcs

Reasons: two phase method also has redistribution step and redistribution of data using range or hash partitioning can easily end up with skew.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



1. **Volume** → Sessions 1, 2, 3, 4
 - How to process Big Data in Volume?

2. **Complexity** → Sessions 5, 6, 7, 8
 - How to apply machine learning algorithms to every aspect of Big Data?

3. **Velocity** → Sessions 9, 10, 11
 - How to handle and process Fast Streaming Data?

<https://tutorcs.com>

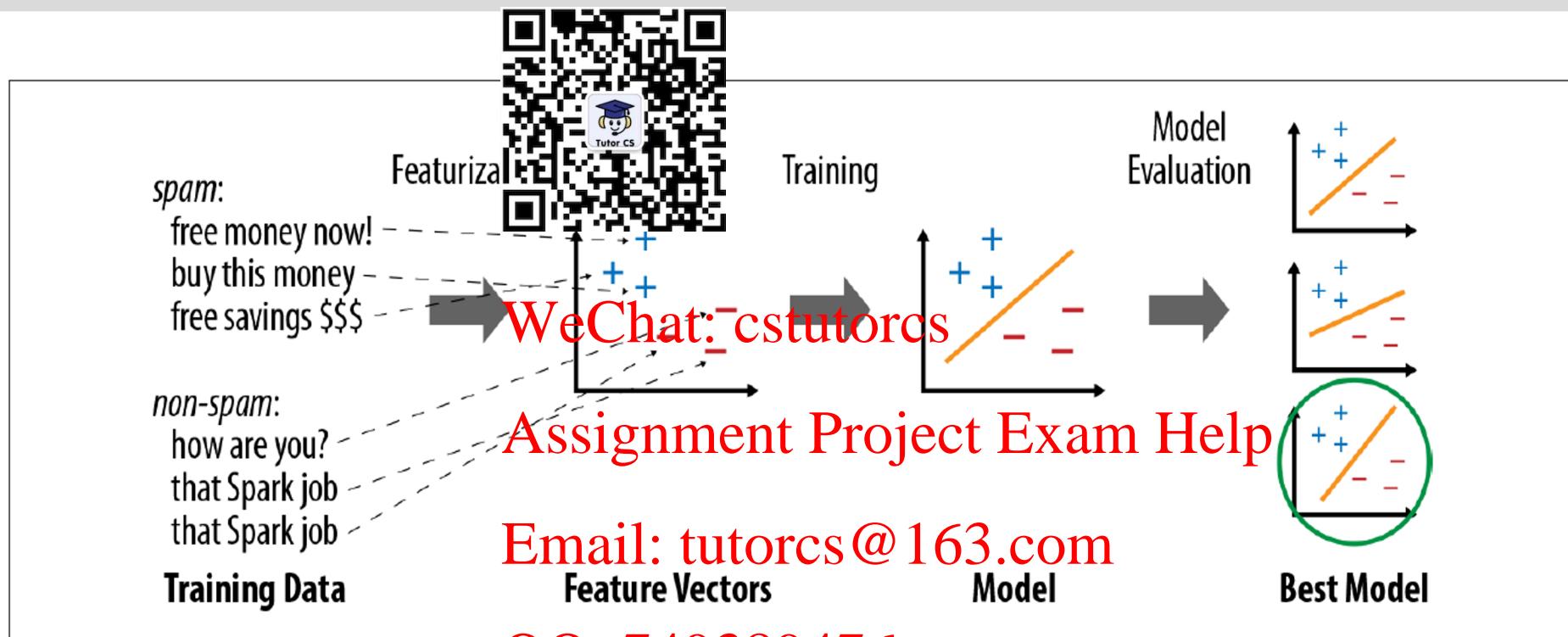


Figure 11-1. Typical steps in a machine learning pipeline

<https://tutorcs.com>

Machine Learning: Featurization

- **Extraction:** Extracting features from “raw” data
 - Count Vectorizer
 - TF-IDF
 - Word2Vec
- **Transformation:** Scaling, converting, or modifying features
 - Tokenization
 - Stop Words Removal
 - String Indexing
 - One Hot Encoding
 - Vector Assembler
- **Selection:** Selecting a subset from a larger set of features
 - Vector Slicer



WeChat: [tutorcs](#)

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: [749389476](#)

<https://tutorcs.com>

Flux Quiz 8

程序代写代做 CS编程辅导

In a product recommendation system, adding another feature (e.g., realizing that which book you should recommend to a user also depends on which movies she's watched) could give a large improvement.



Solution: True

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Types of Machine Learning



- Supervised,
- Unsupervised

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Types of Machine Learning: Supervised

- In supervised machine learning, the data consists of a set of input records.
- Each of these records have **associated labels**.
- The goal is to predict the output label(s) given a new unlabeled input.
- Two types of supervised machine learning:
 1. *Classification* and
 2. *Regression*.



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Supervised Machine Learning: Classification



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476



Binary classification example: dog or not dog

<https://tutorcs.com>

Supervised Machine Learning: Classification



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

QQ: 749389476

Multinomial classification example: Australian shepherd, golden retriever, or poodle

<https://tutorcs.com>

Decision Trees: To Jog or Not To Jog 程序代写代做 CS编程辅导

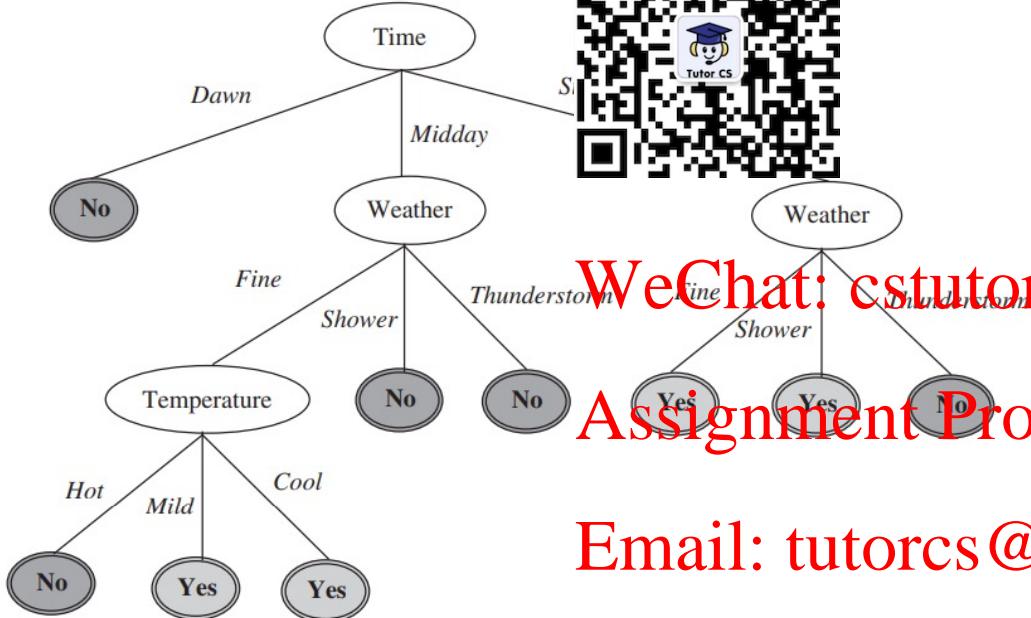


Figure 17.15 Final decision tree

WeChat: cstutorcs
Assignment Project Exam Help

A decision tree is constructed based only on the given training dataset. It is not based on a universal belief.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

ID3 (Iterative Dichotomiser 3)

- It constructs DT, by finding node attribute that returns the highest information gain to split



Steps

- Compute the entropy for dataset S $\rightarrow H(S)$
- For every attribute/feature A :
 - Calculate entropy for each categorical value of A $\rightarrow H(S_i)$
 - Take weighted average entropy for the current attribute $H(S, A) = \sum_{i \in Values(A)} p_i H(S_i)$
 - Calculate IG for the current attribute $\rightarrow IG(S, A) = H(S) - H(S, A)$
- Pick the attribute with highest IG to be a node, and split dataset by its branch to child nodes/subsets $\rightarrow QQ: 749389476$
- Repeat same process at every child node until the tree is complete $\rightarrow https://tutorcs.com$

Stopping condition: when data in each partition have same target class

Flux Quiz 9

程序代写代做 CS编程辅导

What is the **entropy** for the data set in the table below?



Rec#	Weather	Temperature	Time	Day	Jog (target class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

QQ: 749389476

<https://tutorcs.com>

63

data set in the table

Entropy: Measure of uncertainty or randomness in data
Low value -> Less uncertainty, high value -> high uncertainty

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutores@163.com

$$IG(S, A) = H(S) - H(S, A)$$
$$= H(S) - \sum_{i \in Values(A)} p_i H(S_i)$$

Entropy before
(on entire set A)

Entropy after a decision
based on A

S_i Subset/partition of data after splitting S

Weighted
sum entropy
given A

ID3

程序代写代做 CS编程辅导

Entropy for the given probability of the target class

$$\sum_{i=1}^n p_i = 1, \text{ can be calculated as follows:}$$

$$\text{entropy}(p_1, p_2, \dots, p_n) = \sum_{i=1}^n (p_i \log(1/p_i))$$



, ..., p_n where

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

$$\text{entropy}(\text{Yes}, \text{No}) = 5/15 \cdot \log(15/5) + 10/15 \cdot \log(15/10)$$

$$= 0.2764$$

(17.2)
WeChat: cstutorcs

Assignment Project Exam Help

(17.3)
Email: tutorcs@163.com

- Step 1: Calculate entropy for the training dataset in Figure 17.11. The result is previously calculated as 0.2764 (see equation 17.3).

QQ: 749389476

Jog	
Yes	No
5	10

<https://tutorcs.com>

ID3

程序代写代做 CS编程辅导

$$\begin{aligned} \text{entropy}(\text{Weather}=\text{Fine}) &= 4/7 \times \log(7/4) \\ &= 0.2966 \end{aligned}$$



(17.4)

$$\begin{aligned} \text{entropy}(\text{Weather}=\text{Shower}) &= 1/4 \times \log(4/1) \\ &= 0.2442 \end{aligned}$$

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Weather	Jog	
	Yes	No
Fine	4	3
Shower	1	3
Thunderstorm	0	4
		15

程序代写代做 CS编程辅导

Weighted sum entropy (*Weather*) = Weighted entropy(*Fine*) + Weighted entropy(*Shower*) + Weighted entropy(*Thunderstorm*)
 $= 7/15 \times 0.2966 + 4/15 \times 0.2442 + 4/15 \times 0.2442 = 0.2035$



(17.6)

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Rec#	Weather	Temperature	Time	Day	Jog (<i>Target Class</i>)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

	Weather	Jog		15
		Yes	No	
	Fine	4	3	7
	Shower	1	3	4
	Thunderstorm	0	4	4

程序代写代做 CS编程辅导

$$\begin{aligned}
 \text{gain}(\text{Weather}) &= \text{entropy}(\text{training dataset } D) \\
 &= 0.2764 - 0.2035 \\
 &= 0.0729
 \end{aligned}$$



attribute Weather

(17.7)

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

Figure 17.11. Training dataset

程序代写代做 CS编程辅导



- Step 3: Process attribute *Temperature*

- Calculate weighted sum entropy of attribute *Temperature*:

$$\text{entropy}(\text{Hot}) = 2/5 \times \log(5/2) + 3/5 \times \log(5/3) = 0.2923$$

$$\text{entropy}(\text{Mild}) = \text{entropy}(\text{Hot})$$

$$\text{entropy}(\text{Cool}) = 1/5 \times \log(5/1) + 4/5 \times \log(5/4) = 0.2173$$

$$\begin{aligned} \text{weighted sum entropy}(\text{Temperature}) &= 5/15 \times 0.2923 + 5/15 \times 0.2173 \\ &= 0.2674 \end{aligned}$$

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 7.11. Training dataset

	Temperature	Jog		15
		Yes	No	
	Hot	2	3	5
	Mild	3	2	5
	Cool	1	4	5

程序代写代做 CS编程辅导



- Step 4: Process attribute *Time*

- Calculate weighted sum entropy of attribute *Time*:

$$\text{entropy}(\text{Dawn}) = 0 + 5/5 \times \log(5/5) = 0$$

$$\text{entropy}(\text{Midday}) = 2/6 \times \log(6/2) + 4/6 \times \log(6/4) = 0.2764$$

$$\text{entropy}(\text{Sunset}) = 3/4 \times \log(4/3) + 1/4 \times \log(4/1) = 0.2443$$

$$\text{weighted sum entropy } (\text{Time}) = 0 + 6/15 \times 0.2764 + 4/15 \times 0.2443 = 0.1757$$

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

- Calculate information gain for attribute *Time*:

$$\text{gain } (\text{Temperature}) = 0.2764 - 0.1757 = 0.1007$$

QQ: 749389476

<https://tutorcs.com>

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 7.11. Training dataset

	Time	Jog		15
		Yes	No	
	Dawn	0	5	5
	Midday	2	4	6
	Sunset	3	1	4

程序代写代做 CS编程辅导



Step 5: Process attribute *Day*

- Calculate weighted sum entropy of attribute *Day*:

$$\begin{aligned} \text{entropy(Weekday)} &= 4/10 \times \log(10/4) + 6/10 \times \log(10/6) \\ &= 0.2923 \end{aligned}$$

$$\begin{aligned} \text{entropy(Weekend)} &= 1/5 \times \log(5/1) + 4/5 \times \log(5/4) \\ &= 0.2173 \end{aligned}$$

$$\begin{aligned} \text{weighted sum entropy (Day)} &= 10/15 \times 0.2923 + 5/15 \\ &\quad \times 0.2173 = 0.2674 \end{aligned}$$

- Calculate information gain for attribute *Day*:

$$\text{gain}(\text{Temperature}) = 0.2764 - 0.2674 = 0.009$$

QQ: 749389476

<https://tutorcs.com>

WeChat: cstutorcs

Email: tutorcs@163.com

Assignment Project Exam Help

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 7.11. Training dataset

Day	Jog		Yes	No	Count
	Weekend	Weekday			
Weekend	4	6	10	15	15
Weekday	1	4	5		

程序代写代做 CS编程辅导

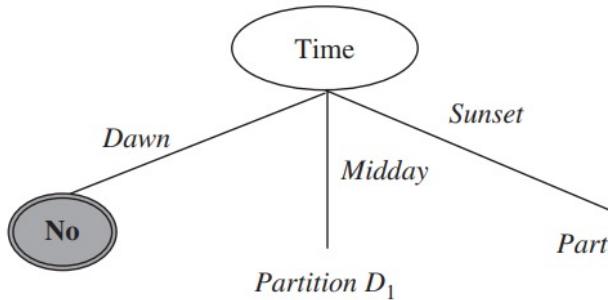


Figure 17.13. Attribute *Time* as the root node

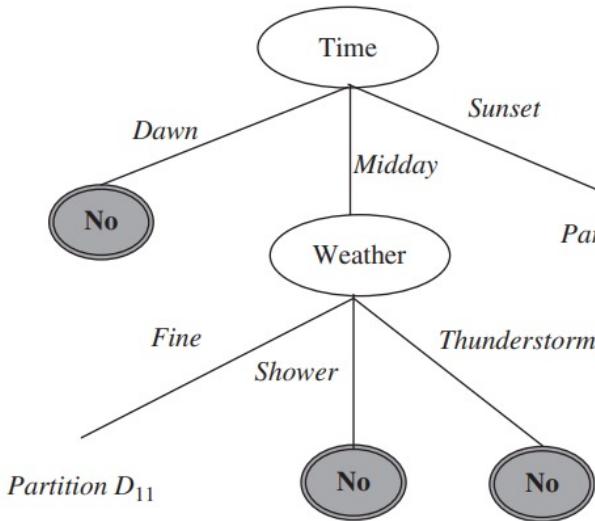
Assignment Project Exam Help
WeChat: cstutorcs
Email: tutorcs@163.com
QQ: 749389476

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Comparing equations 17.7, 17.8, 17.9, and 17.10 ,and 17.10 for the gain of each other attributes (Weather, Temperature, Time, and Day), the biggest gain is *Time*, with gain value = 0.1007 (see equation 17.9), and as a result, attribute *Time* is chosen as the first splitting attribute. A partial decision tree with the root node *Time* is shown in Figure 17.13.

<https://tutorcs.com>



WeChat: cstutorcs

Assignment Project Exam Help

Figure 17.14 Attribute
Weather as next splitting attribute

- The next stage is to process partition D_1 consisting of records with Time=Midday. Training dataset partition D_1 consists of 6 records with record#s: 3, 6, 8, 9, 10, and 15. The next task is to determine the splitting attribute for partition D_1 , whether it is Weather, Temperature, or Day.

<https://tutorcs.com>

72

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Decision Trees: To Jog or Not To Jog 程序代写代做 CS编程辅导

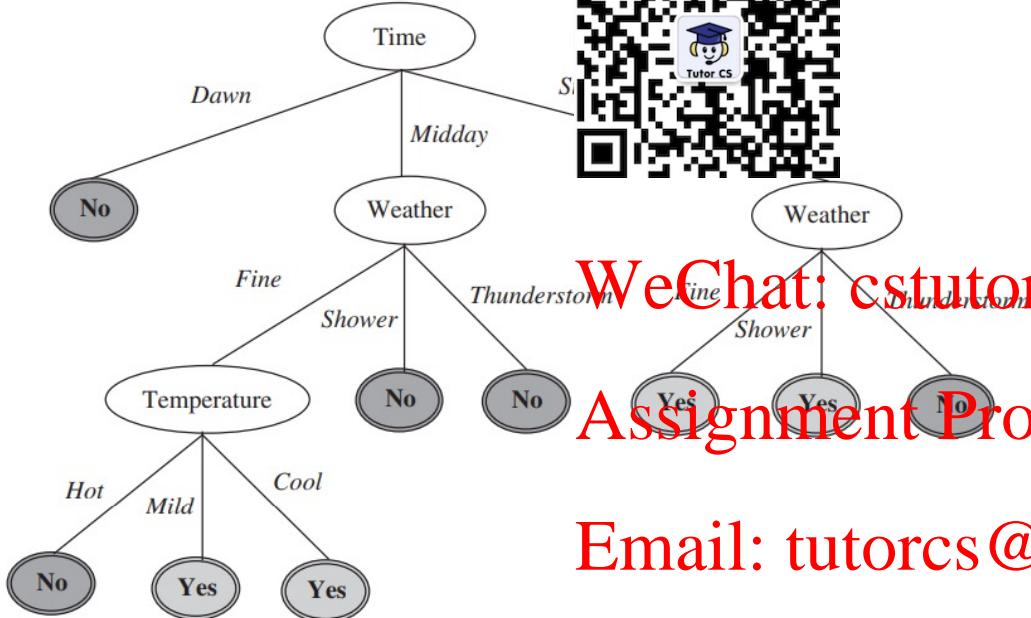


Figure 17.15 Final decision tree

WeChat: cstutorcs
Assignment Project Exam Help

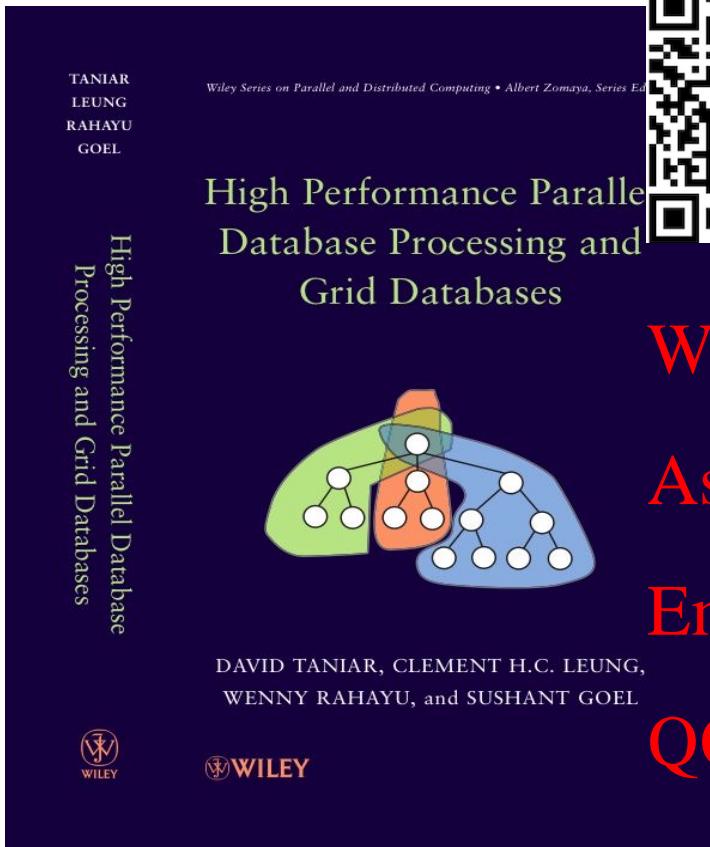
A decision tree is constructed based only on the given training dataset. It is not based on a universal belief.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Complexity → Session 程序代写, 代做 CS 编程辅导
5, 6, 7,



Chapter 17 Parallel Clustering and Classification

WeChat: cstutorcs

Assignment Project Exam Help

17.1 Clustering and Classification

17.2 Parallel Clustering

17.3 Parallel Classification

17.4 Summary

17.5 Biographical Notes

17.6 Exercises

QQ: 749389476

<https://tutorcs.com>

Parallel Classification: Decision Tree

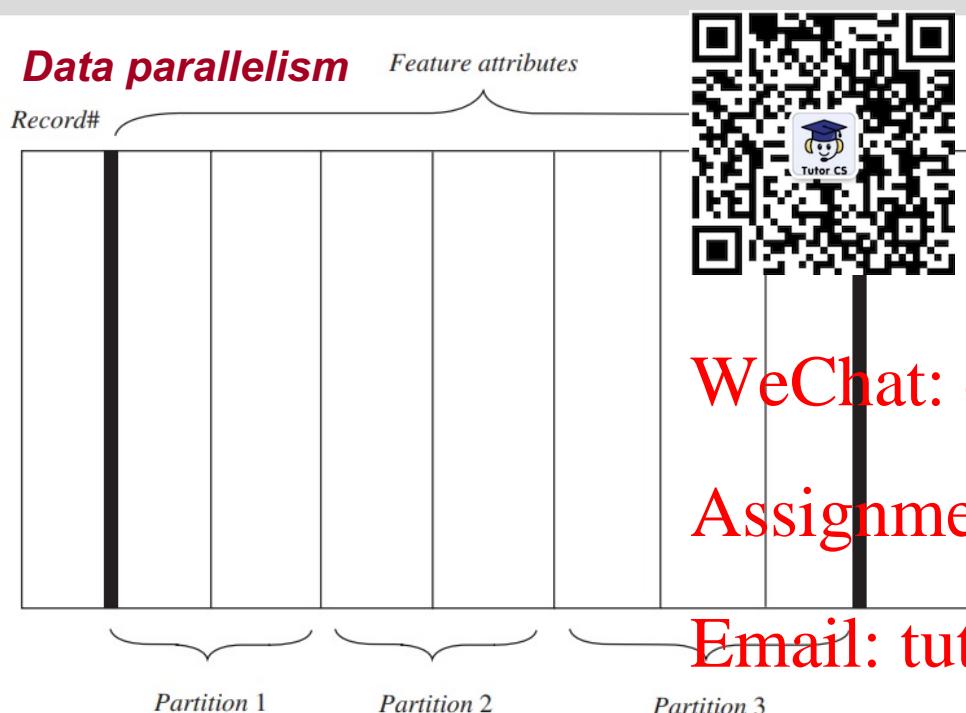


Figure 17.16 Vertical data partitioning of training data set
QQ: 749389476

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Parallel Classification: Decision Tree

Data parallelism:

Rec#	Weather	Temperature	Class
1	Fine	Mild	No
2	Fine	Hot	No
3	Shower	Mild	No
4	Thunderstorm	Cool	No
5	Shower	Hot	Yes
6	Fine	Hot	No
7	Fine	Cool	No
8	Thunderstorm	Cool	No
9	Fine	Cool	Yes
10	Fine	Mild	Yes
11	Shower	Hot	No
12	Shower	Mild	No
13	Fine	Cool	No
14	Thunderstorm	Mild	No
15	Thunderstorm	Hot	No



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

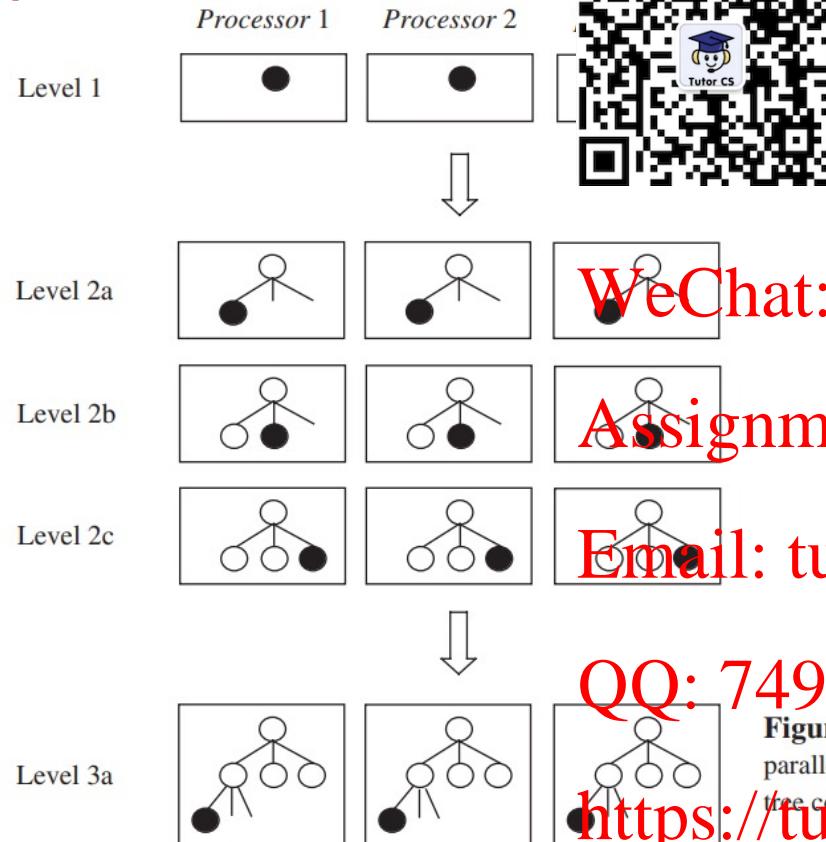
Partition 1

<https://tutorcs.com>

Partition 2

Parallel Classification: Decision Tree

Data parallelism



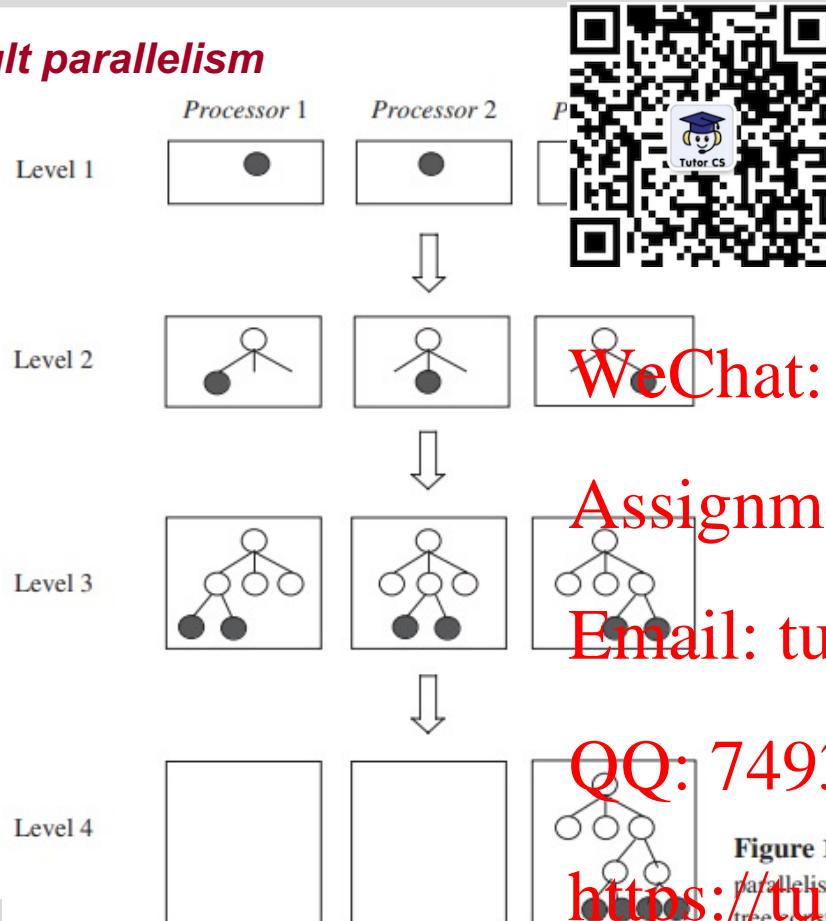
Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

- All processors will process one particular node at each level or sub-level at a time
- Each processor will compute IGs for certain attributes (in parallel), which are then shared to determine splitting
- ‘Intra tree node parallelism’

Parallel Classification: Decision Tree

Result parallelism



Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

❑ Multiple nodes are processed concurrently using several processors at each level

❑ Main rule: A processor that processes a child node will also its parent nodes

❑ ‘Inter tree node parallelism’

Types of Machine Learning: Unsupervised

- Instead of predicting, unsupervised ML helps you to better understand the structure of your data.
- Two types of unsupervised machine learning:
 1. *Clustering* and
 2. *Association*.



WeChat: cstutorcs

Assignment Project Exam Help

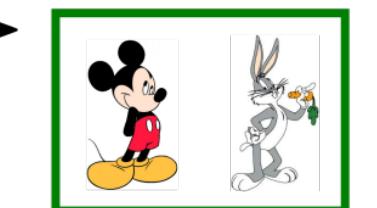
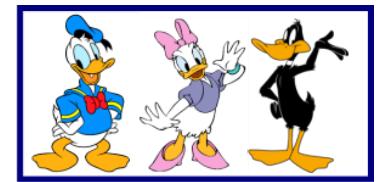
Email: tutores@163.com

QQ: 749389476

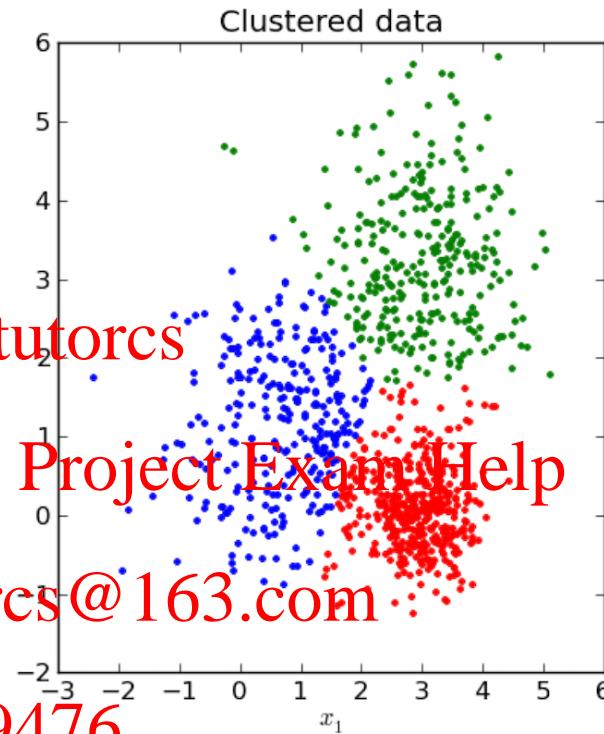
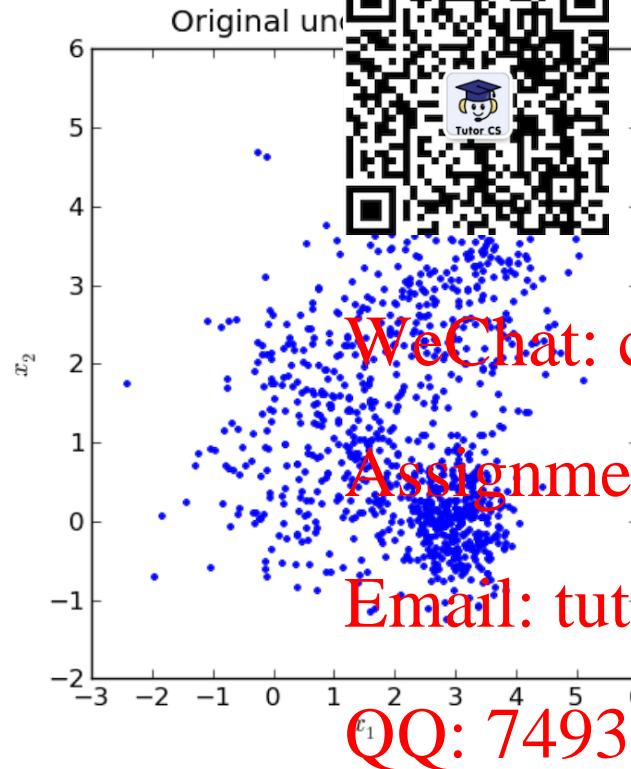
<https://tutorcs.com>



Unsupervised
Learning



Unsupervised Machine Learning: Clustering



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

Clustering example

<https://tutorcs.com>

Algorithm k-Means:

- (**Initialization**) Specified number of clusters, and guesses the k seed cluster centers
- (**Assignment Step**) Assign each data point to the cluster with the closest centroid
- Current clusters may receive or lose their members
- (**Update Step**) Each cluster must re-calculate the mean (centroid)
- The process is repeated until the clusters are stable (no change of members)



Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

Algorithm: k-means

Input:

$D = \{x_1, x_2, \dots, x_n\}$

//Data objects

k //Number of desired clusters

Output:

K //Set of clusters

1. Assign initial values for means m_1, m_2, \dots, m_k
2. Repeat
3. Assign each data object x_i to the cluster which has the closest mean
4. Calculate new mean for each cluster
5. Until convergence criteria is met

Flux Quiz 10

程序代写代做 CS编程辅导

Data D = {5, 19, 25, 21, 4, 1, 7, 6, 10, 2, 20, 14, 11, 27, 9, 3, 16}

Number of clusters: k = 3

Initial centroids: m₁=6, m₂=8.



Which of the following grouping is correct after applying K-Means algorithm?

WeChat: cstutorcs

Assignment Project Exam Help

Solution:

$$C_1 = \{1, 2, 3, 4, 5, 6\}$$

Email: tutorcs@163.com

$$C_2 = \{7, 8, 9, 10, 11, 14\}$$

$$C_3 = \{16, 17, 19, 20, 21, 23, 25, 27\}$$

QQ: 749389476

<https://tutorcs.com>

k-Means: Step-By-Step Example



- Data $D = \{5, 19, 25, 11, 4, 17, 23, 8, 7, 6, 10, 2, 20, 14, 11, 27, 9, 3, 16\}$

- Number of clusters $K=3$

- Initial centroids: $m_1=6, m_2=7$, and $m_3=8$

- **First Iteration**

- Clusters:

- $C_1=\{1, 2, 3, 4, 5, 6\}$

- $C_2=\{7\}$

- $C_3=\{8, 9, 10, 11, 14, 16, 17, 19, 20, 21, 23, 25, 27\}$

- Re-calculated centroids: $m_1=3.5, m_2=7$, and $m_3=16.9$

First Iteration: Calculating euclidean distance, determining the cluster membership and calculating new centroid.

D	5	19	25	11	4	17	23	8	7	6	10	2	20	14	11	27	9	3	16	
d(m1, Di)	1	13	19	15	2	5	11	17	2	1	0	4	4	14	8	5	21	3	3	10
d(m2, Di)	2	12	18	14	3	6	10	16	1	0	1	3	5	13	7	4	20	2	4	9
d(m3, Di)	3	11	17	13	4	7	9	15	0	1	2	2	6	12	6	3	19	1	5	8

<https://tutorcs.com>



MONASH
University

k-Means: Step-By-Step Example



- Clusters:
 - $C_1 = \{1, 2, 3, 4, 5\}$
 - $C_2 = \{7\}$
 - $C_3 = \{8, 9, 10, 11, 14, 16, 17, 19, 20, 21, 23, 25, 27\}$
- New centroids: $m_1 = 3.5$, $m_2 = 7$, and $m_3 = 16.9$
- Second Iteration
 - Clusters:
 - $C_1 = \{1, 2, 3, 4, 5\}$
 - $C_2 = \{6, 7, 8, 9, 10, 11\}$
 - $C_3 = \{14, 16, 17, 19, 20, 21, 23, 25, 27\}$
 - Re-calculated centroids: $m_1 = 3$, $m_2 = 8.5$, and $m_3 = 20.2$

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Second Iteration: Calculating euclidean distance, determining the cluster membership and calculating new centroid.



MONASH UNIVERSITY	5	19	25	21	4	1	17	23	8	7	6	10	2	20	14	11	27	9	3	16
d(m1, Di)	1.5	15.5	21.5	17.5	0.5	2.5	13.5	19.5	4.8	3.5	2.5	6.5	1.5	16.5	10.5	7.5	23.5	5.5	0.5	12.5
d(m2, Di)	2	12	18	14	3	6	10	16	1	0	1	3	5	13	7	4	20	2	4	9
d(m3, Di)	11.9	2.1	8.1	4.1	12.9	15.9	0.1	6.1	8.9	9.9	10.9	6.9	14.9	3.1	2.9	5.9	10.1	7.9	13.9	0.9

k-Means: Step-By-Step Example



- Clusters:
 - $C_1 = \{1, 2, 3, 4, 5\}$
 - $C_2 = \{6, 7, 8, 9, 10, 11, 14\}$
 - $C_3 = \{16, 17, 19, 20, 21, 23, 25, 27\}$
- New centroids: $m_1=3$, $m_2=8.5$, and $m_3=20.2$

Third Iteration

Assignment Project Exam Help

- Clusters:
 - $C_1 = \{1, 2, 3, 4, 5\}$
 - $C_2 = \{6, 7, 8, 9, 10, 11, 14\}$
 - $C_3 = \{16, 17, 19, 20, 21, 23, 25, 27\}$
- Re-calculated centroids: $m_1=3$, $m_2=9.29$, and $m_3=21$

Third Iteration: Calculating euclidean distance, determining the cluster membership and calculating new centroid.

<https://tutorcs.com>

D	M	O	N	A	H	19	25	21	4	1	17	23	8	7	6	10	2	20	14	11	27	9	3	16
d(m1, Di)	2	16	22	18	1	2	14	20	0.5	4	3	7	1	17	11	8	24	6	0	13				
d(m2, Di)	10.5	16.5	12.5	4.5	7.5	8.5	14.5	0.5	1.5	2.5	1.5	6.5	11.5	5.5	2.5	18.5	0.5	5.5	7.5					
d(m3, Di)	15.2	1.2	4.8	0.8	16.2	19.2	3.2	2.8	12.2	13.2	14.2	10.2	18.2	0.2	6.2	9.2	6.8	11.2	17.2	4.2				

k-Means: Step-By-Step Example



- Clusters:
 - $C_1 = \{1, 2, 3, 4, 5, 6\}$
 - $C_2 = \{6, 7, 8, 9, 10, 11, 12, 13, 14\}$
 - $C_3 = \{16, 17, 19, 20, 21, 23, 25, 27\}$
- New centroids: $m_1=3$, $m_2=9.29$, and $m_3=21$
- **Fourth Iteration** Assignment Project Exam Help
 - Clusters:
 - $C_1 = \{1, 2, 3, 4, 5, 6\}$
 - $C_2 = \{7, 8, 9, 10, 11, 14\}$
 - $C_3 = \{16, 17, 19, 20, 21, 23, 25, 27\}$
 - Re-calculated centroids: $m_1=3.5$, $m_2=9.83$, and $m_3=21$

Fourth Iteration: Calculating euclidean distance, determining the cluster membership and calculating new centroid.

D	5	19	25	21	4	1	17	23	8	7	6	10	2	20	14	11	27	9	3	16
d(m1, Di)	2	16	22	18	1	2	14	20	5	Ω	4	3	7	1	17	11	8	24	6	13
d(m2, Di)	9.7	15.7	11.7	5.3	8.3	7.7	13.7	1.3	2.3	3.3	0.7	7.3	10.7	4.7	1.7	17.7	0.3	6.3	6.7	
d(m3, Di)	16.0	2.0	4.0	0.0	17.0	20.0	4.0	2.0	13.0	14.0	15.0	11.0	19.0	1.0	7.0	10.0	6.0	12.0	18.0	5.0

k-Means: Step-By-Step Example



- Clusters:
 - $C_1 = \{1, 2, 3, 4, 5, 6\}$
 - $C_2 = \{7, 8, 9, 10, 11, 12, 13, 14\}$
 - $C_3 = \{16, 17, 19, 20, 21, 23, 25, 27\}$
- New centroids: $m_1 = 6$, $m_2 = 9.83$, and $m_3 = 21$
- **Fifth Iteration**
 - No data movement from clusters (Process Terminated)

Assignment Project Exam Help

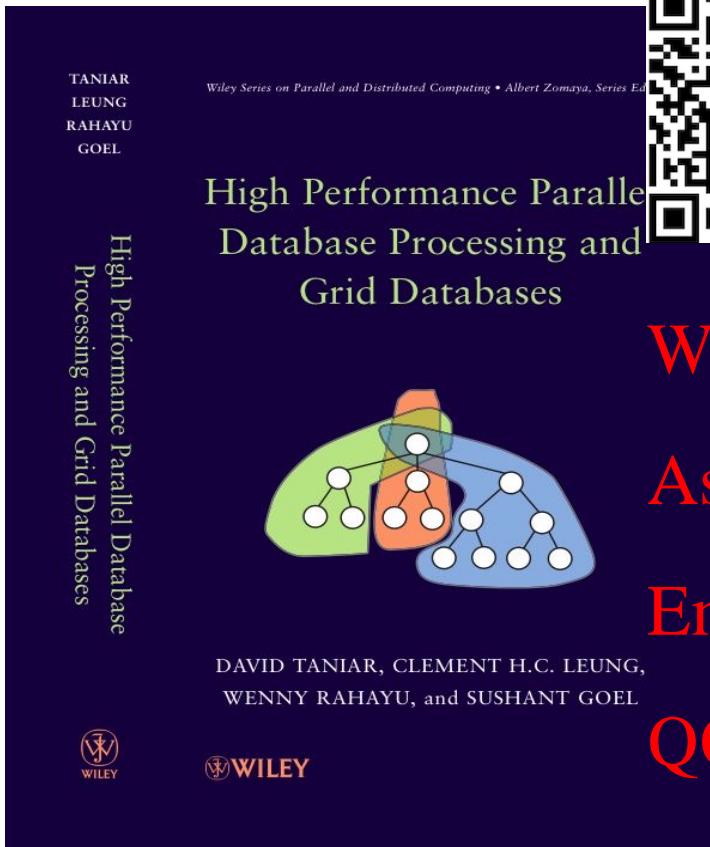
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

m_1	m_2	m_3	C_1	C_2	C_3
6	7	8	1, 2, 3, 4, 5, 6	7	8, 9, 10, 11, 14, 16, 17, 19, 20, 23, 25, 27
3.5	7	16.9	1, 2, 3, 4, 5	6, 7, 8, 9, 10, 11	14, 16, 17, 19, 20, 21, 23, 25, 27
3	8.5	20.2	1, 2, 3, 4, 5	6, 7, 8, 9, 10, 11, 14	16, 17, 19, 20, 21, 23, 25, 27
3	9.29	21	1, 2, 3, 4, 5, 6	7, 8, 9, 10, 11, 14	16, 17, 19, 20, 21, 23, 25, 27
3.5	9.83	21	1, 2, 3, 4, 5, 6	7, 8, 9, 10, 11, 14	16, 17, 19, 20, 21, 23, 25, 27

Complexity → Session 程序代写, 代做 CS 编程辅导
5, 6, 7,



Chapter 17 Parallel Clustering and Classification

WeChat: cstutorcs

Assignment Project Exam Help

17.1 Clustering and Classification

17.2 Parallel Clustering

17.3 Parallel Classification

17.4 Summary

17.5 Biographical Notes

17.6 Exercises

QQ: 749389476

<https://tutorcs.com>

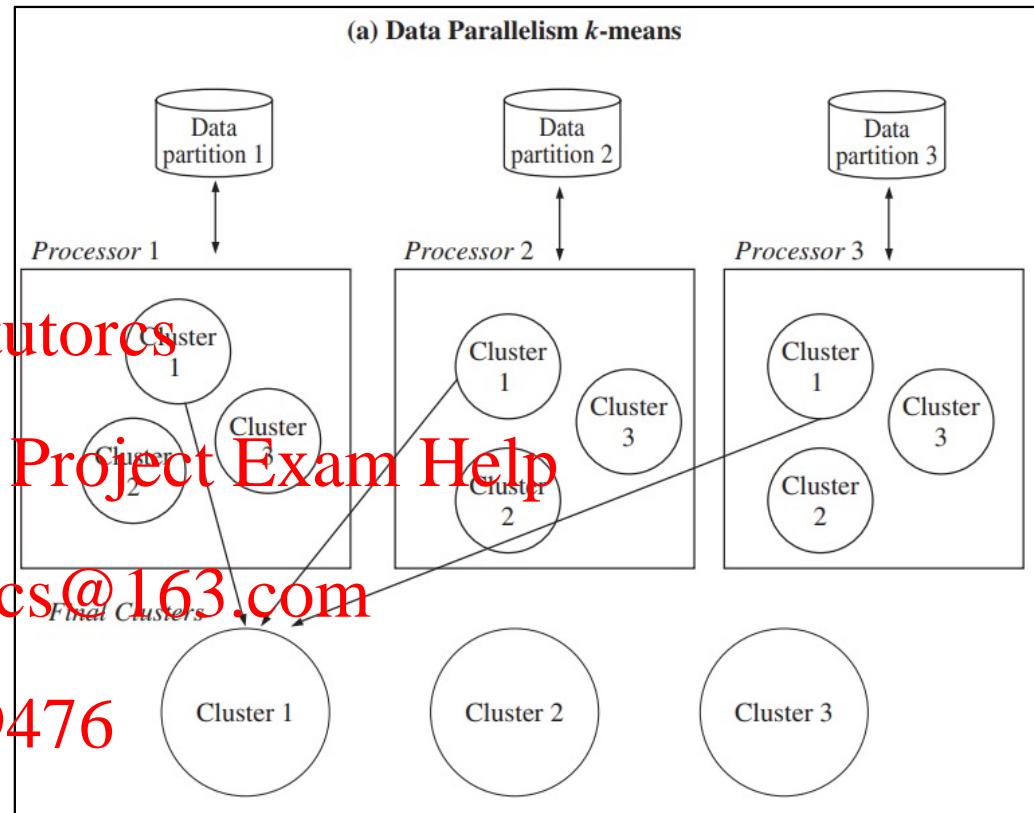
Parallel K-means clustering

Data parallelism of k-means

- Data is partitioned into multiple partition
- Each processor will work independently to create three clusters
- The final clusters from each processor are respectively united



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



Parallel K-means clustering



- **Result Parallelism** of
- Focuses on clusters partitioning
- Each processor will work on a particular target cluster
- During the iteration, the memberships of cluster can change. -→ data movement across processors

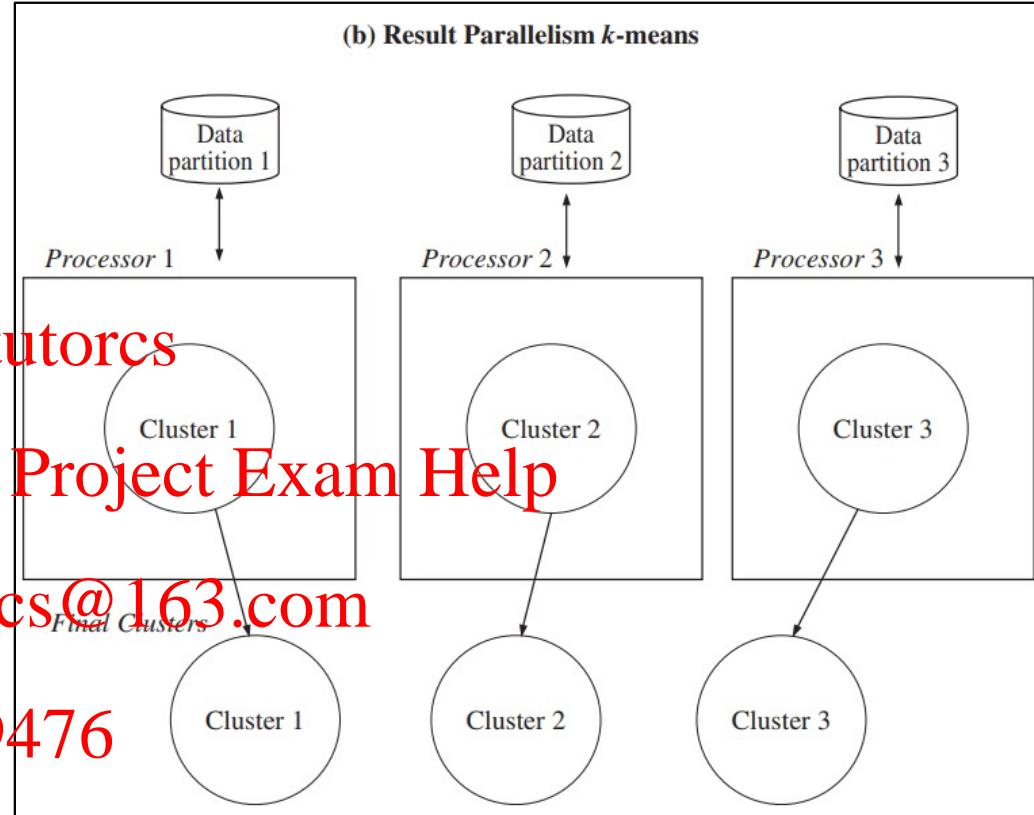
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Collaborative Filtering – User based

It calculates the similarity between users to make implicit recommendation.



Steps:

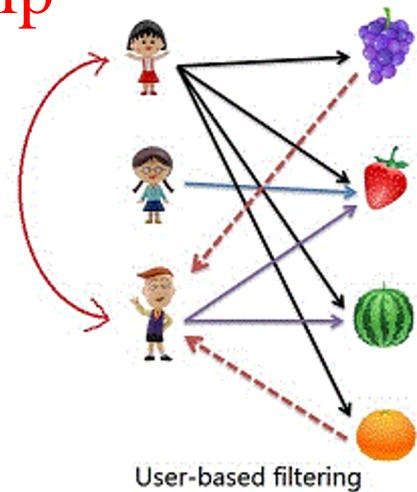
1. Calculate the similarity between PersonA and all other users.
2. Predict the ratings of items for PersonA based on similar users.
3. Select top-N rated items for PersonA.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Similar
users



Flux Quiz 11

程序代写代做 CS编程辅导

Collaboration Filtering:  Example

Name	Star Trek	Superman	Batman	Hulk
Harry	4	2	?	5
John	5	3	4	?
Rob	3	?	4	4

Aim: Recommend top-2 movies
to Harry

QQ: 749389476

<https://tutorcs.com>

Collaboration Filtering: Walkthrough Example (user-based)

Step 1: Calculate the similarity between Harry and all other users



Name	Star Trek	Superman	Batman	Hulk
Harry	4	2	?	5
John	5	3	4	?
Rob	3	?	4	4

Cosine similarity

Email: tutorcs@163.com

$$sim(u, u') = \cos(\theta) = \frac{\mathbf{r}_u \cdot \mathbf{r}_{u'}}{\|\mathbf{r}_u\| \|\mathbf{r}_{u'}\|} = \sum_i \frac{r_{ui} r_{u'i}}{\sqrt{\sum_i r_{ui}^2} \sqrt{\sum_i r_{u'i}^2}}$$

r_{ui} - value of rating user u gives to item i

<https://tutorcs.com>

Collaboration Filtering: Walkthrough Example (user-based)

Step 1: Calculate the similarity between Harry and all other users

Name	Star Trek	Star wars	Supernatural	Hulk
Harry	4	2	?	5
John	5	3	4	3
Rob	3	?	4	3



WeChat: cstutorcs

Assignment Project Exam Help

Cosine similarity

$$\text{Sim}(\text{Harry}, \text{John}) = \frac{(4*5)+(2*3)+(4*3)}{\sqrt{4^2+2^2+4^2} * \sqrt{5^2+3^2+3^2}} = 0.97$$

Email: tutorcs@163.com

QQ: 749389476

$$\begin{aligned}\text{Sim}(\text{Harry}, \text{Rob}) &= \frac{(4*3)+(5*4)+(4*3)}{\sqrt{4^2+5^2+4^2} * \sqrt{3^2+4^2+3^2}} \\ &= 1.00\end{aligned}$$

<https://tutorcs.com>

Collaboration Filtering: Walkthrough Example (user-based)

Step 2: Predict the ratings of user Harry



Name	Star Trek	Star wars	Superman	Hulk
Harry	4	2	?	5
John	5	3	4	?
Rob	3	?	4	4

WeChat: cstutorcs

Predicted rating is calculated based on aggregation of some similar users' rating of the item

$$r_{u,i} = k \sum_{u' \in U} \text{simil}(u, u') r_{u',i}$$

with normalising factor

Assignment Project Exam Help
 $k = 1 / \sum_{u' \in U} |\text{simil}(u, u')|$,

Calculate k as a normalising factor $k = \frac{1}{(0.97+1)} = 0.51$

$$\begin{aligned} R(\text{Harry}, \text{Superman}) &= k * ((\text{sim}(\text{Harry}, \text{John}) * R(\text{John}, \text{Superman})) + (\text{sim}(\text{Harry}, \text{Rob}) * R(\text{Rob}, \text{Superman}))) \\ &= 0.51((0.97 * 4) + (1 * 4)) = 4.02 \end{aligned}$$

<https://tutorcs.com>

Collaboration Filtering: Walkthrough Example (user-based)

Step 3: Select top-2 rated movies for Harry

Name	Star Trek	Avatar	Superman	Batman	Hulk
Harry	4	5	4.02	5	4
John	5	3	4	?	3
Rob	3	?	4	4	3

WeChat: cstutorcs

$\text{Top-2}(Harry, \text{movies}) = \text{Batman, Superman}$

Email: tutorcs@163.com

QQ: 749389476

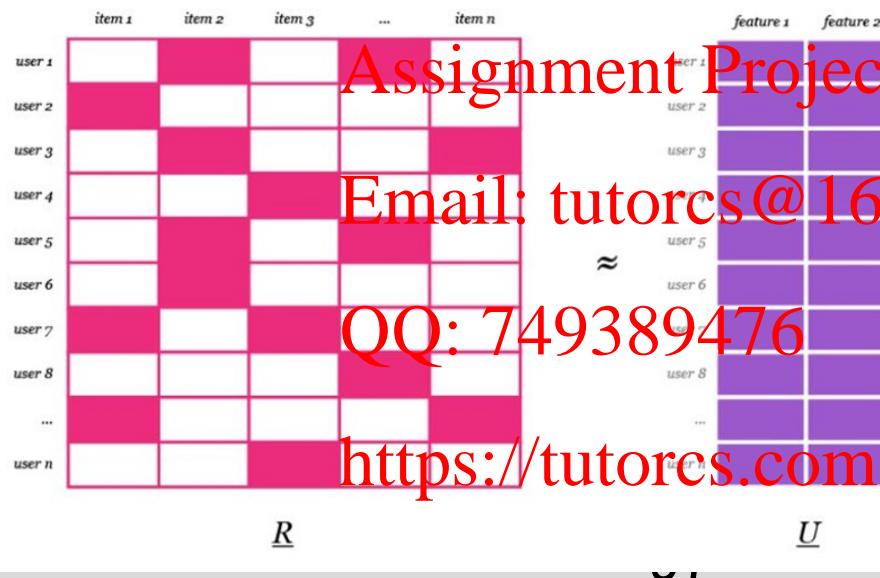
<https://tutorcs.com>

Model-based Collaborative Filtering



- Latent factor model based CF: Extract (latent) user and item profiles through **matrix factorization**
- **Matrix factorization:** Factor a rating matrix into some smaller representation of the original matrix through alternating least squares. The product of lower dimensional matrices equals the original one
- Example: Factor the rating matrix \mathbf{R} into user matrix \mathbf{U} and item matrix \mathbf{V}

WeChat: cstutorcs



Alternating least squares

- ALS method aims to estimate the factor matrices (\mathbf{U} & \mathbf{V}) such that it will approximate the original rating matrix.
- This is achieved by minimizing the root mean square error (RMSE) between the original ratings and the predicted values.



WeChat: cstutorcs



Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

User factor matrix, \mathbf{U}
<https://tutorcs.com>

98
<https://developers.google.com/machine-learning/recommendation/collaborative/matrix>

ALS Procedure:

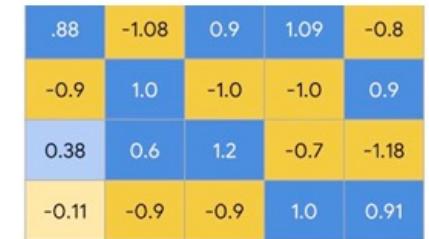
Optimizing alternately to find \mathbf{U} , \mathbf{V}

- Randomly initialize \mathbf{U} and \mathbf{V}
- Iterating the following steps:
 - Fixing \mathbf{U} → Optimizing \mathbf{V}
 - Fixing \mathbf{V} → Optimizing \mathbf{U}
- Each iteration will minimize the square errors

Item factor matrix, \mathbf{V}



≈



Predicted rating matrix $\hat{\mathbf{R}}$



1. **Volume** → Sessions 1, 2, 3, 4
 - How to process Big Data in Volume?

2. **Complexity** → Sessions 5, 6, 7, 8
 - How to apply machine learning algorithms to every aspect of Big Data?

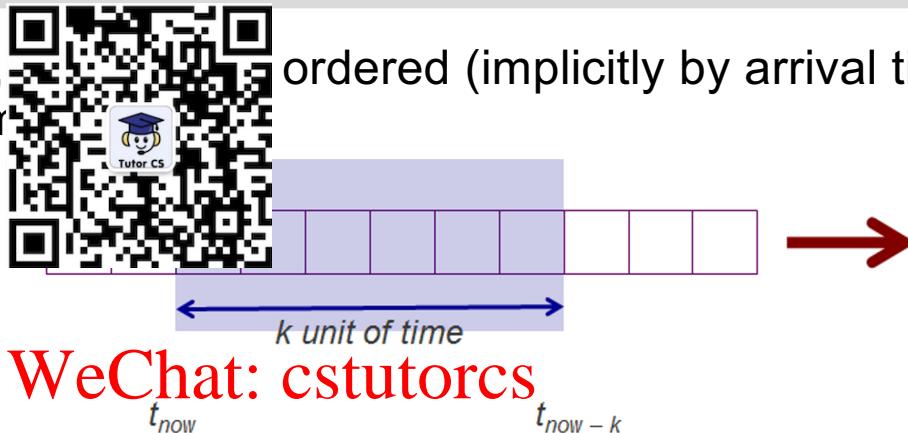
3. **Velocity** → Sessions 9, 10, 11
 - How to handle and process Fast Streaming Data?

<https://tutorcs.com>

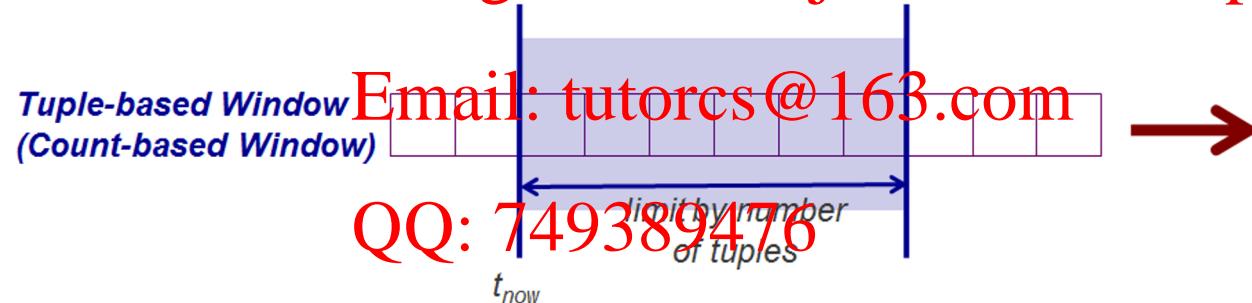
程序代写代做 CS 编程辅导

Windowing System in Unbounded Streams

A **data stream** is a real-time, ordered (implicitly by arrival time or explicitly by timestamp) sequence of items.



Assignment Project Exam Help



<https://tutorcs.com>

Stream Window – Time Based Window



WeChat: cstutorcs Size of window is specified by fixed time duration

Overlapping Sliding Window

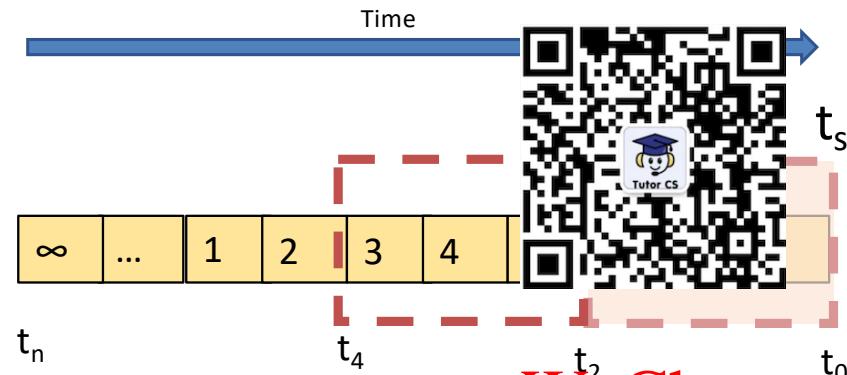
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Stream Window – Time Based Window



- Window size is based on time, eg 2 seconds
- Window can be advanced by:
 - $t_e - t_s$ (window size)
 - Duration less than the window size (sliding window).
- In uniform data rate, the number of tuples will be the same for each window.

Non-Overlapping Sliding Window

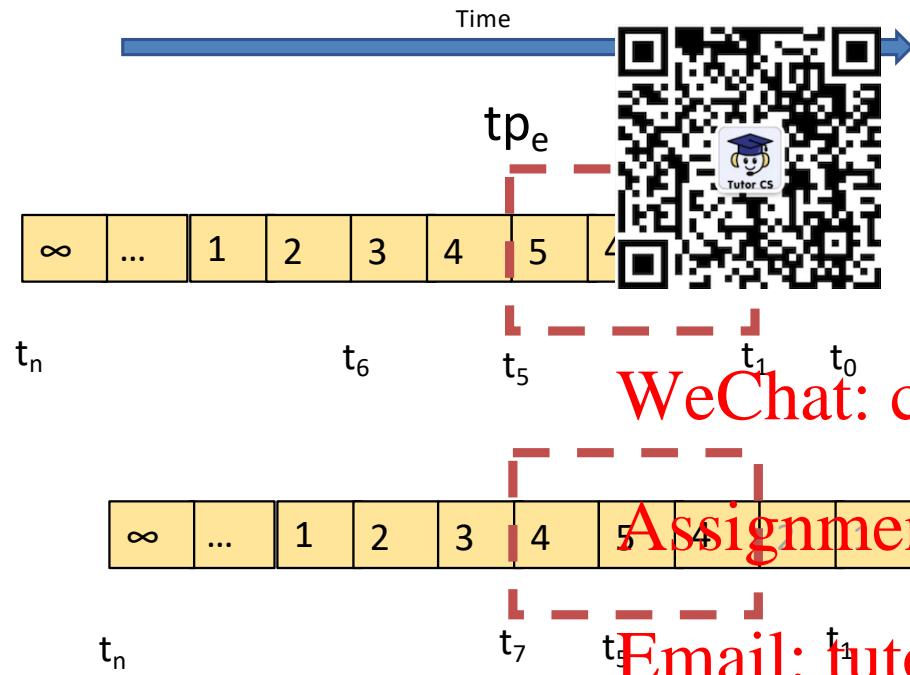
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Stream Window – Tuple Based Window



Window size is 3 tuples
Slide window by 1 tuple

WeChat: cstutorcs Size of window is specified by number of tuples in the window

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Database vs Stream Processing

- Bounded data (assessable).
- Relatively static data
- Complex, ad-hoc query
- Possible to backtrack during processing
- Exact answer to a query
- Tuples arrival rate is low



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

- Unbounded data (data keeps coming without end).
- Dynamic data.
- Simple, continuous query
- No backtracking, single pass operation.
- Approximate answer to a query
- Tuples arrival rates is high

Sliding Window: Produce an approximate answer to a data stream query is to evaluate the query not over the entire past history of data streams, but rather only over sliding windows of recent data from the streams.

Event vs Processing time



- Event time: the time when data is produced by the source.
- Processing time: the time when data is arrived at the processing server.
- In ideal situation, event time equals processing time.
- In real world event time is earlier than the processing time due to network delay.
- The delay can be uniform (ideal situation) or non-uniform (most of real network situation).
- Data may arrive in “burst” (bursty network)

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Velocity → Session 9, 10, 11 程序代写代做 CS 编程辅导

Joins In Data Streams

- Symmetric Hash Join
- M Join
- AM Join
- Handshake Join



WeChat: cstutorcs

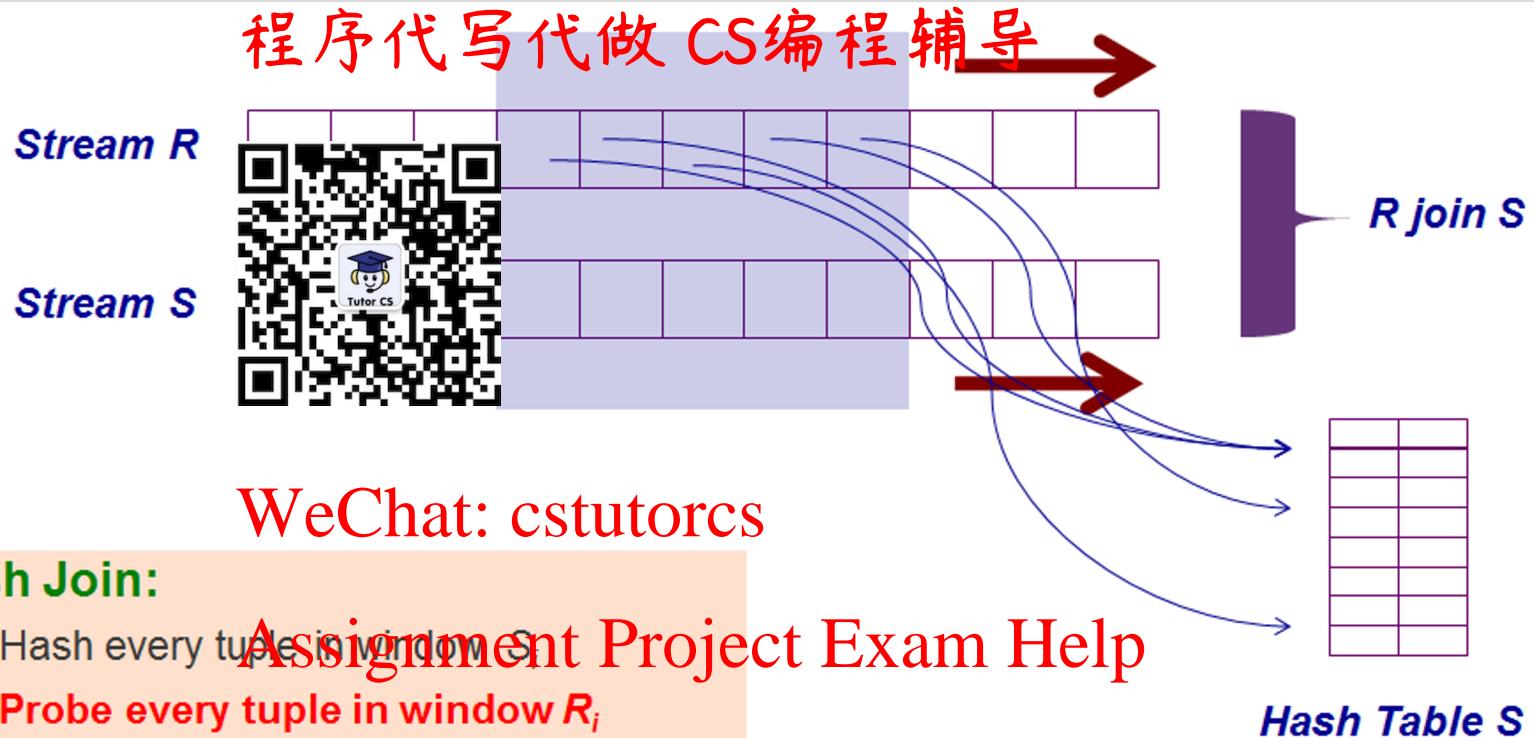
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Hash Join



QQ: 749389476

Problem: If *r* comes in first and then *s* comes in later, then *r* will not be compared with the *s*
→ If *r* and *s* has same key, we will miss join operation

<https://tutorcs.com>

Symmetric Hash Join

Step 1:

tuple r arrives

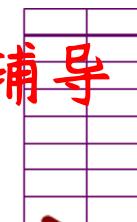


Stream R

tuples s not yet
processed

Stream S

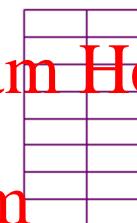
WeChat: estutorcs



Hash Table R



R join S



Hash Table S

Assignment Project Exam Help

Symmetric Hash Join Process:

Email: tutorcs@163.com

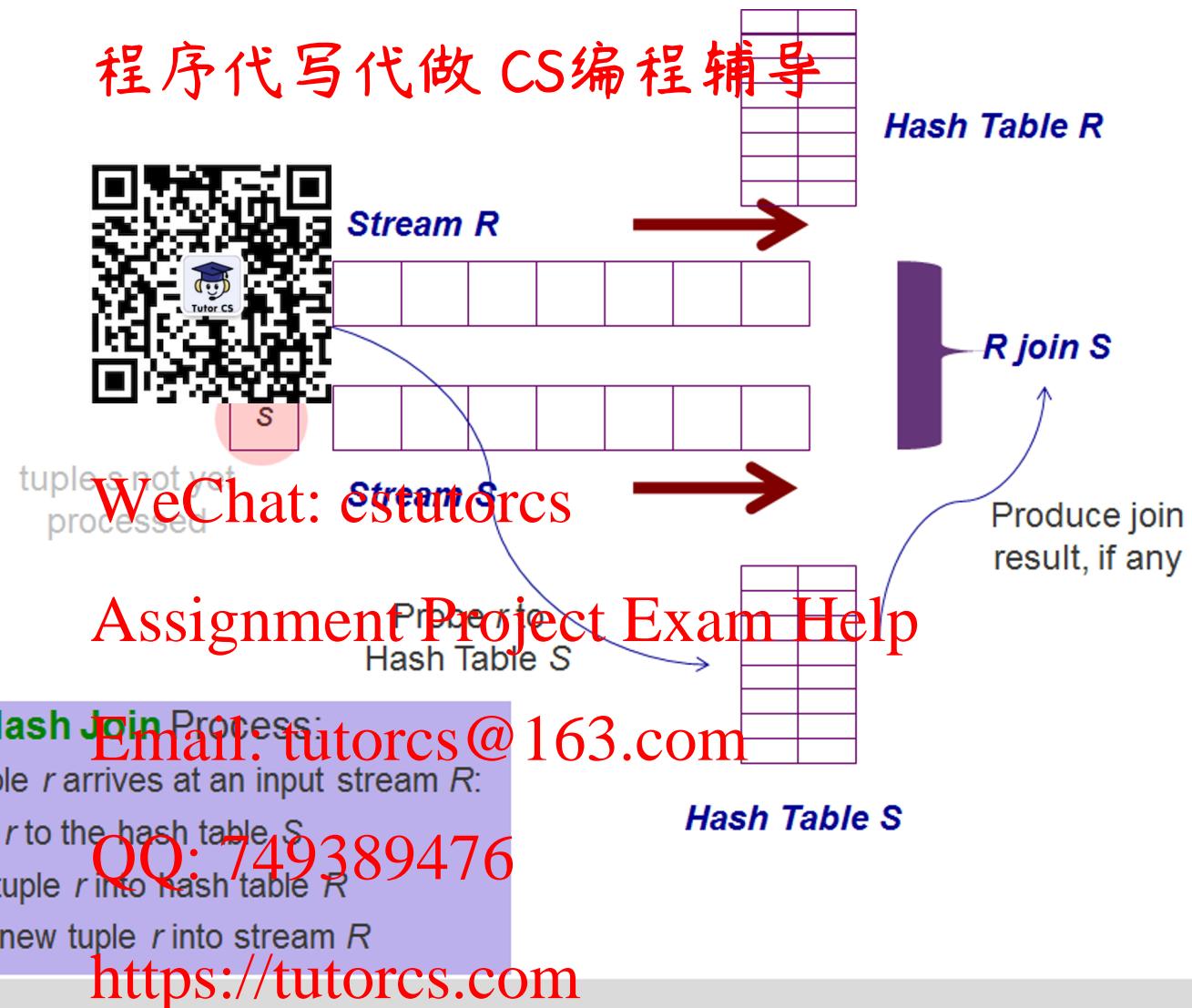
- When a tuple r arrives at an input stream R :
 - Probe r to the hash table S
 - Hash tuple r into hash table R
 - Insert new tuple r into stream R

QQ: 749389476

<https://tutorcs.com>

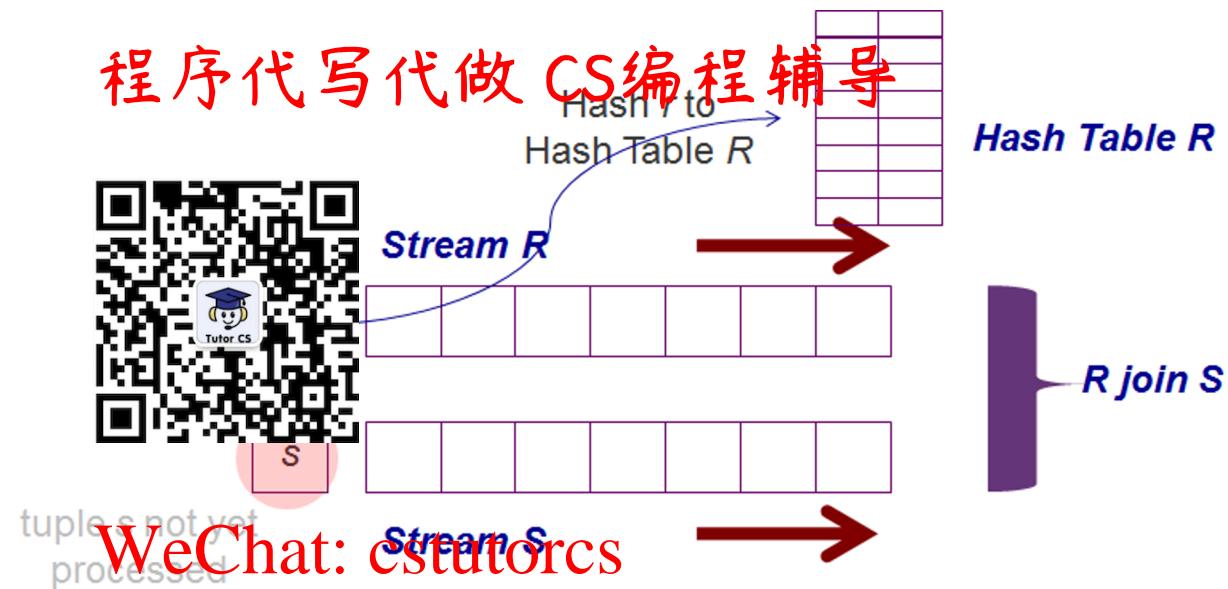
Symmetric Hash Join

Step 2:



Symmetric Hash Join

Step 3:



Assignment Project Exam Help

Symmetric Hash Join Process:
Email: tutorcs@163.com

- When a tuple r arrives at an input stream R :
 - Probe r to the hash table S
 - Hash tuple r into hash table R
 - Insert new tuple r into stream R

QQ: 749389476

<https://tutorcs.com>

Hash Table S

Symmetric Hash Join

Step 4:



Assignment Project Exam Help

Symmetric Hash Join Process:

Email: tutorcs@163.com

- When a tuple r arrives at an input stream R :
 - Probe r to the hash table S
 - Hash tuple r into hash table R
 - Insert new tuple r into stream R

QQ: 749389476

<https://tutorcs.com>

M-Join

程序代写代做 CS编程辅导

Step 1:

tuple r arrives



Stream R



Hash Table R

tuples s and t not
yet processed

Stream S

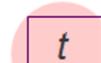




Hash Table S

Assignment Project Exam Help

Stream T





Hash Table T

M-Join Process:

- When a tuple r arrives at an input stream R :
 - **Probe r to all other hash tables**
 - Hash tuple r into hash table R
 - Insert new tuple r into stream R

QQ: 749389476

<https://tutorcs.com>

M-Join

Step 2:

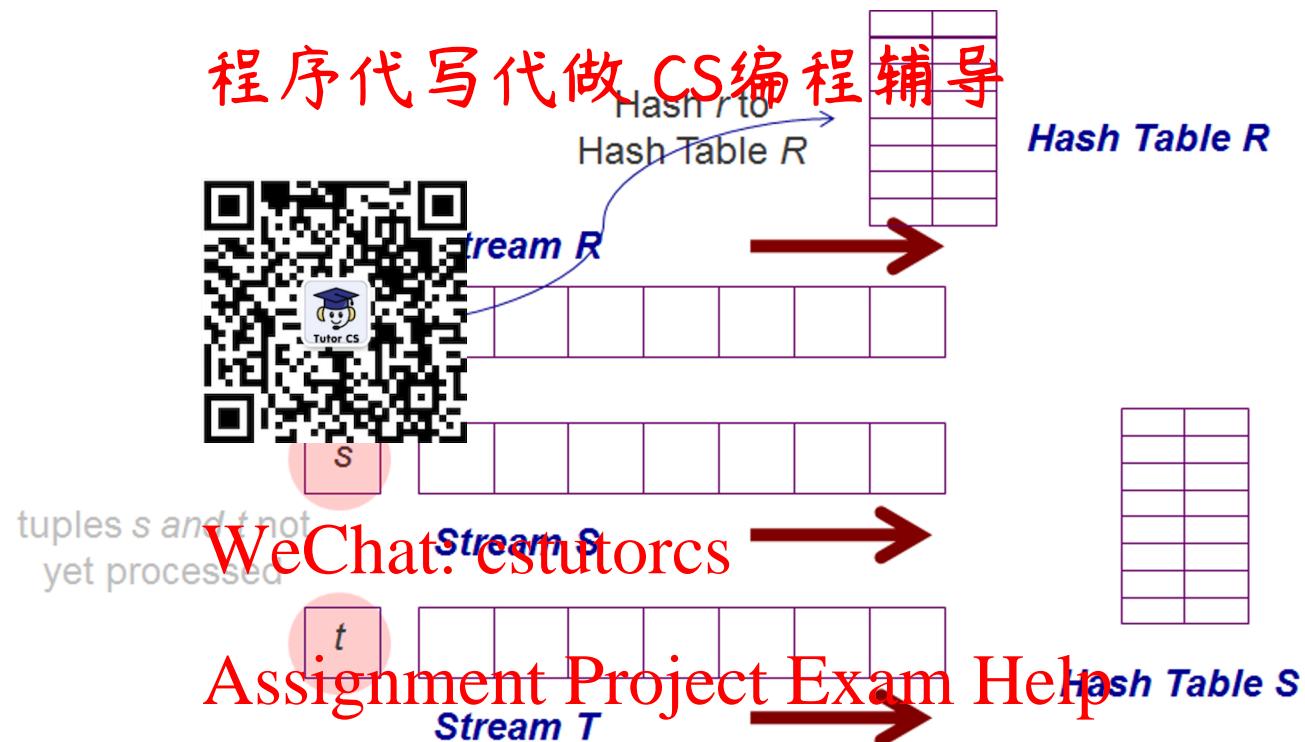


M-Join Process:

- When a tuple r arrives at an input stream R :
 - Probe r to all other hash tables
 - Hash tuple r into hash table R
 - Insert new tuple r into stream R

M-Join

Step 3:



M-Join Process:

- When a tuple r arrives at an input stream R :
 - **Probe r to all other hash tables**
 - Hash tuple r into hash table R
 - Insert new tuple r into stream R

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

M-Join

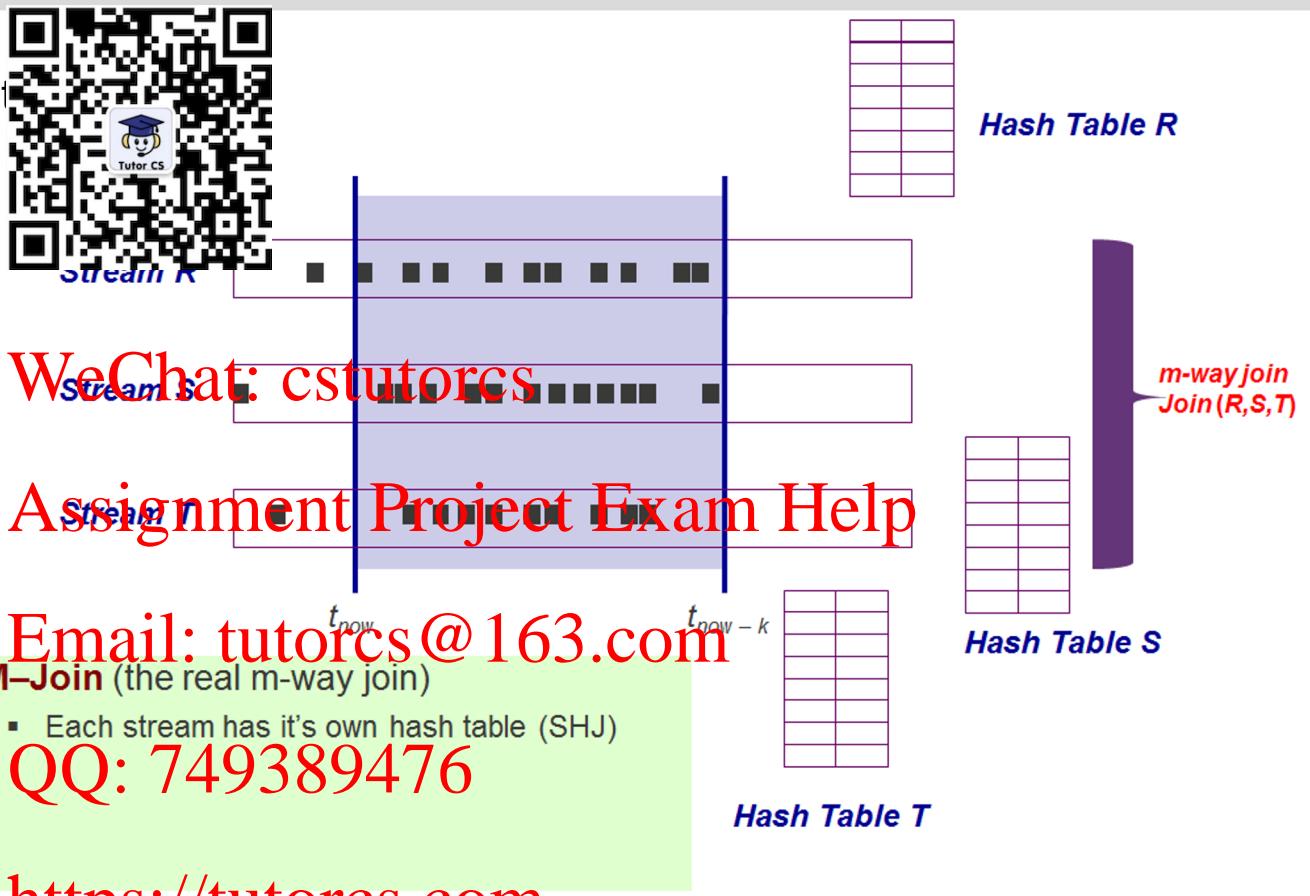
Step 3:



Tuple Slide (Using M-Join)

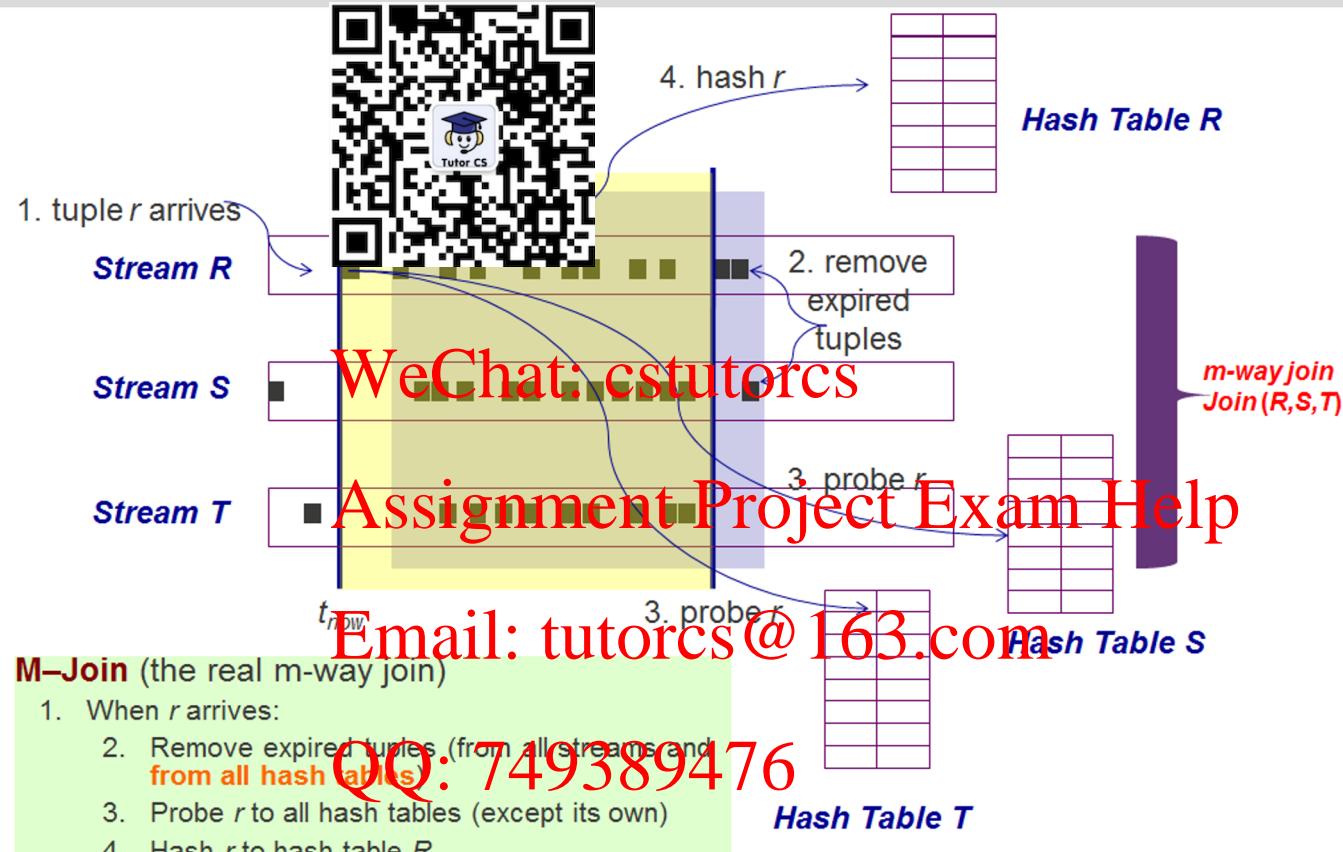
How to join more than two streams at same time?

- Use **M-Join** – a multiway streaming join
- Based on symmetric hash join (work similarly to two-stream)



Tuple Slide (Using M-Join)

程序代写代做 CS编程辅导



Time Slide (Using M-Join) 程序代写代做 CS编程辅导



How to slide the window?

- Tuple Slide
- Time Slide

QQ: 749389476

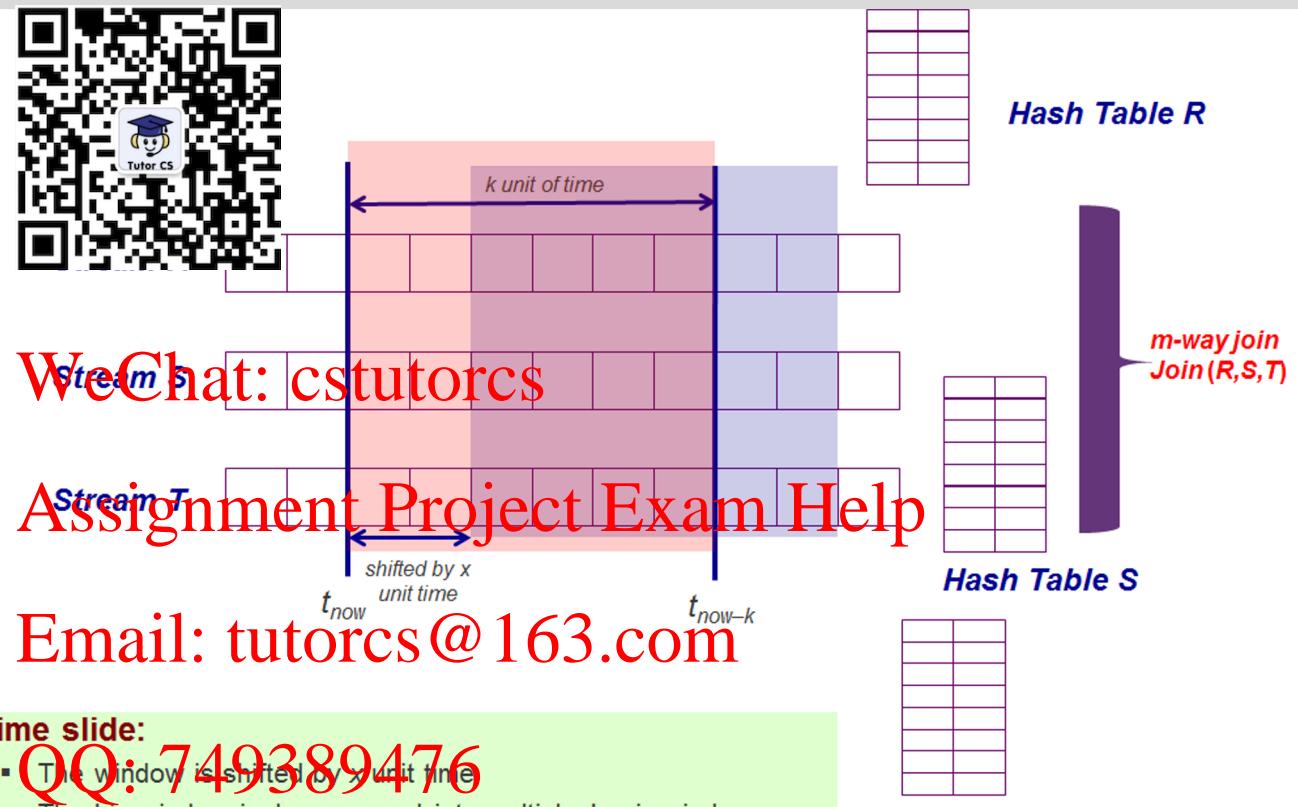
<https://tutorcs.com>

Time Slide (Using M-Join)

程序代写代做 CS编程辅导

- Let one box here represents one unit time or one basic window.

Ex: 1 unit time = 1 minute
- Window size = 6 mins
- slide = 2mins



Time slide:

- The window is shifted by x unit time
- The big window is decomposed into multiple basic windows
- The big window is then shifted by 1 or more basic windows

<https://tutorcs.com>

Time Slide (Using M-Join) 程序代写代做 CS编程辅导



Time Slide (using M-Join):

- We use the current window (assume the previous expired basic windows have been removed as part of the previous join process)
- In the above example, the latest 2 basic windows are new basic windows
- Note that the hash tables do not yet have tuples from these 2 basic windows

QQ: 749389476

<https://tutorcs.com>

Time Slide (Using M-Join) 程序代写代做 CS编程辅导



Time Slide (Using M-Join) 程序代写代做 CS编程辅导



Time Slide (using M-Join):

- **Process Stream R:**
 - Take all tuples from basic windows r_1 and r_2 , and probe to hash tables S and T
 - Take all tuples from basic windows r_1 and r_2 to hash table R

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

Time Slide (Using M-Join) 程序代写代做 CS编程辅导



Time Slide (Using M-Join) 程序代写代做 CS编程辅导



Time Slide (using M-Join):

• Process Stream S:

- Take all tuples from basic windows s_1 and s_2 , and probe to hash tables *R* and *T*
- Take all tuples from basic windows s_1 and s_2 to hash table *S*

QQ: 749389476

<https://tutorcs.com>

Time Slide (Using M-Join)

程序代写代做 CS编程辅导



Time Slide (Using M-Join) 程序代写代做 CS编程辅导



Handshake Join

程序代写代做 CS编程辅导



- After a handshake of a pair of tuples is done, both streams R and S move forward at the same time
- As a result, we will miss a handshake because one tuple from R moves to the right, and one tuple from S moves to the left, and one handshake will be missed

Handshake Join

程序代写代做 CS编程辅导

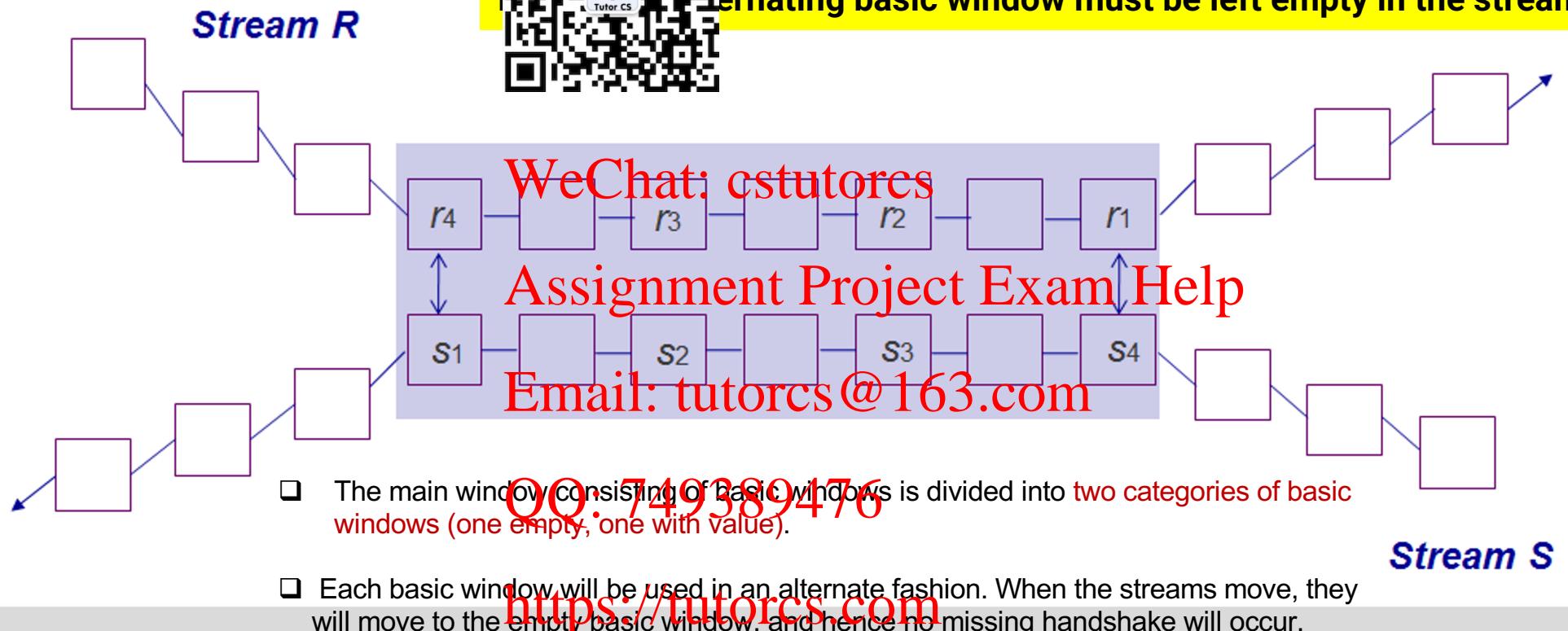


Handshake Join (Solution 1)



Solution 1:

Alternating basic window must be left empty in the stream.



Handshake Join (Solution 1)



Solution 1:



程序代写代做 CS编程辅导

Handshake Join (Solution 1)



程序代写代做 CS编程辅导

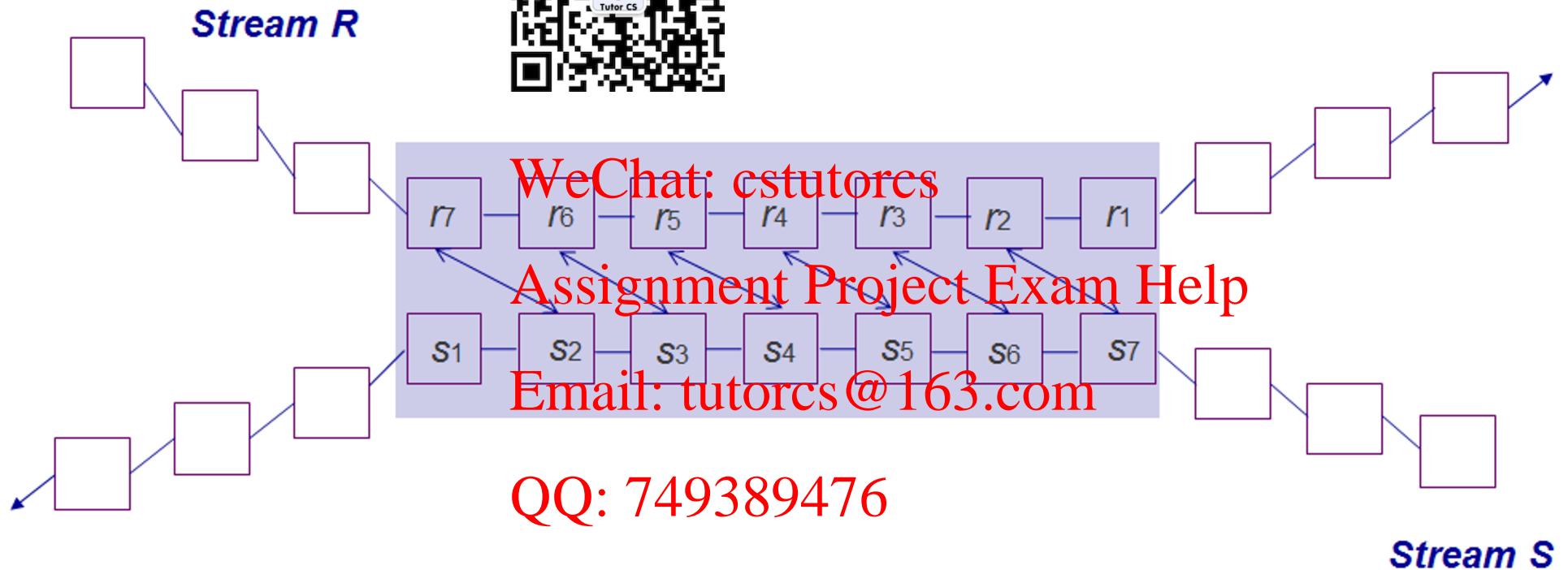
Handshake Join (Solution 2)



Main Idea: After a handshake happens, the streams do not move forward, but each tuple in the stream performs a handshake with the next tuple, and then after that both streams move forward.
<https://tutorcs.com>

程序代写代做 CS编程辅导

Handshake Join (Solution 2)



程序代写代做 CS编程辅导

Handshake Join (Solution 2)



Solution 2:



Velocity → Session 9, 10, 11 程序代写代做 CS 编程辅导



Granularity Reduction In Data Streams

- Group By and Aggregation

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Granularity

程序代写代做 CS编程辅导

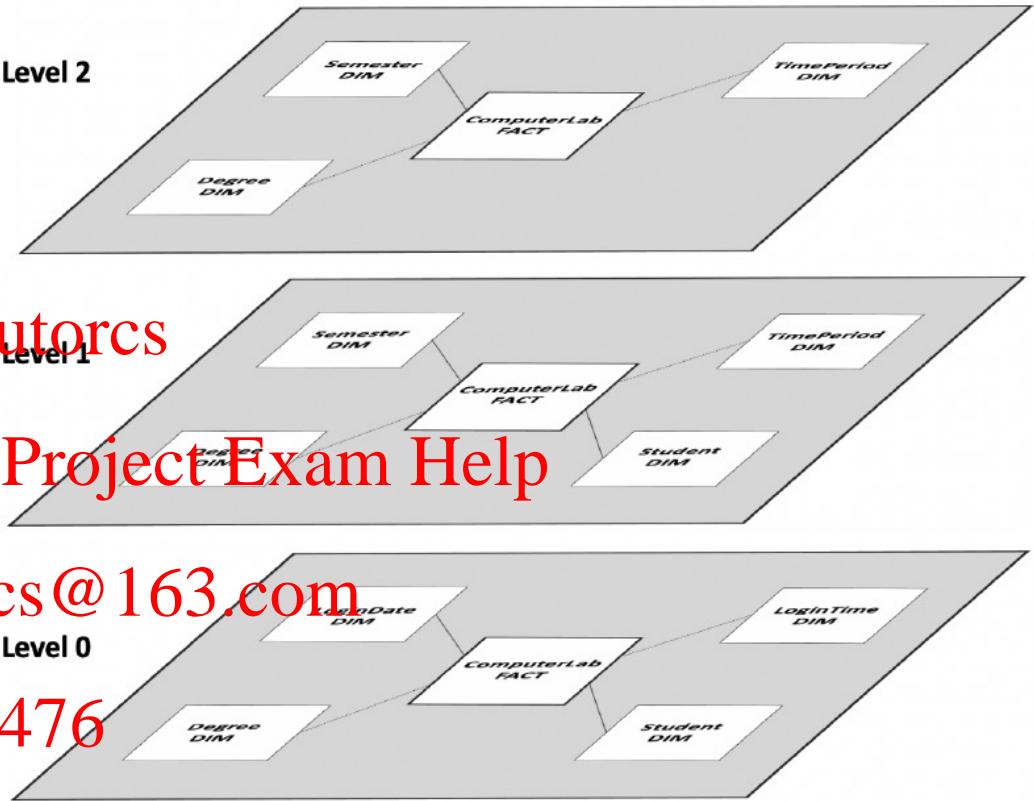


- **Granularity** is the level at which **data** are stored in a database.
- level-0, the bottom level indicating no aggregation.
- level-1 and level-2 with more aggregation.

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com
QQ: 749389476

<https://tutorcs.com>



Flux Quiz 13

程序代写代做 CS编程辅导

In a sales scenario, if one record contains a month sales amount, a window of 6 months is used to calculate the running average sales amount. In this case, the window size is 6 records, and the slide is effective. The number of records in the moving average will be the same as the original records. If one year has 12 records of sales, the moving average will also contain 12 records. Hence, no reduction in terms of the number of records. This is a pure moving average (also known as rolling mean).

The above-mentioned case is an example of:

- A. Overlapped Windows - No granularity reduction
- B. Overlapped Windows - With granularity reduction
- C. Non-Overlapped Windows

<https://tutorcs.com>

Mixed Levels of Granularity

程序代写代做 CS编程辅导

- Different levels of granularities combined into one level.
- Mixed level of granularity can be two types:



- Temporal-based Mixed Levels of Granularity
 - Time based.
- Spatial-based Mixed Levels of Granularity
 - Space or location based.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Sensor Arrays

程序代写代做 CS编程辅导



- A sensor array is a group of sensors, usually deployed in a certain geometry pattern.
- A network of distributed sensors.
- They add new dimension to the observation, and hence it helps to estimate more parameters, to have better picture of the environment being observed, and improve accuracy.
- Two categories:
 1. Multiple sensors measuring the same things, and
 2. Multiple sensors measuring different things, but they are grouped together.

<https://tutorcs.com>

Sensor Arrays

程序代写代做 CS编程辅导



- Multiple sensors measure the same things
- Two methods to lower redundancy of sensor arrays that measure the same thing:
 - Method 1: Reduce and then Merge
 - Method 2: Merge and then Reduce

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

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



- **Multiple sensors measuring different things**
 - Sensors arrays can be a collection of sensors measuring different things within the same environment.
- Example: A simple indoor sensor array, containing three sensors: air quality, temperature, and humidity.
- **Two methods to lower the granularity of sensor arrays that measure the different thing:**
 - Method 1: Reduce, Normalize, and then Merge
 - Method 2: Normalize, Merge and then Reduce

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



1. **Volume** → Sessions 1, 2, 3, 4
 - How to process Big Data of Large Volume?

2. **Complexity** → Sessions 5, 6, 7, 8
 - How to apply machine learning algorithms to every aspect of Big Data?

3. **Velocity** → Sessions 9, 10, 11
 - How to handle and process Fast Streaming Data?

<https://tutorcs.com>

Thank You

程序代写代做 CS编程辅导



Questions?



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>