

02-july-23

M22AI608

Kadapala Rakesh Reddy

ALGORITHMS FOR BIG DATA - CSL7630

MINOR-II

IIT-J. D.C.S.-M.Tech.- Summer sem-I

11

- (i) False: streaming algorithms process data in a memory-efficient manner, often by utilizing limited portions of input at a time.
- (ii) False: streaming algorithms can be used & handles for both Real time data streams & static datasets.
- (iii) True: streaming algorithms can handle parallel processing of data.
- (iv) False: S.A often provide approximate results for computations on streaming data.
- (v) False: universal hashing doesn't mandate the use of cryptographic hash functions.
- (vi) False: universal hashing minimizes collisions but doesn't guarantee unique locations for all keys.
- (vii) False: universal hashing doesn't guarantee constant time complexity for hash table operations.
- (viii) False: Property Testing applies to both Quantitative & Qualitative properties.
- (ix) True: Property Testing algorithms can handle dynamic datasets that change over time.
- (x) False: Property testing aims to provide probabilistic results in sublinear time, not exact solutions.

(2)

Item	2	3	1	4	2	3	1	3	1	4	5	3	2	1	6
Count	1	1	1	1	2	2	2	3	3	2	1	4	3	4	1

Here algorithm initializes a table with a row for each possible item in stream & column for count of no. of times that item has been seen.

The algorithm then iterates through data stream, updating the counts for each item as it goes. If count for an item reaches k , then item is added to output set. If count for an item falls below k , then item is evicted from output set.

The algorithm terminates when the end of data stream is reached. The final output of algorithm is set of items with a count at least k .

So final output set = {1, 2, 3, 4, 5}

The misra-Gries algorithm has several benefits over traditional frequency counting approaches for large stream of data. First it is very memory efficient

- Also, one needs to store a small table with a row for each possible item in stream.
- This is contrast to traditional frequency counting approaches, which may need to store count for each item in stream.

Reservoir sampling 2 Initialize reservoir size S

- For each item in the data stream
- if item not in reservoir
- Remove a random item from a reservoir
- final Q/P is set of items in reservoir

final $Q/P = \{2, 3, 1, 4, 6\}$

Reservoir sampling is a simple & efficient algorithm for randomly selecting a subset of size 'k' from a large stream of data.

The Algo. works by maintaining a reservoir of size 'k'.

- As each item is added to reservoir with probability $\frac{1}{k}$ if the reservoir is full, random item is removed to make a new item.

③ (11) A universal one-way hash function family is a set of hash functions such as that for any 2 distinct keys, the probability that they hash to same value is at most $\frac{1}{m}$ where $m = \text{size of hash table}$.

(12) A solution that utilizes a universal hash functions family to speed up searching a keywords across a large set of text documents could be a hash table. The hash would be initialized with set of hash functions from universal hash function family. When a keyword is searched for, it would be hashed using each of hash functions in family. The results of hash functions would then be used to lookup the keyword in hash table. If the keyword found in hash table, then it is returned, else keyword is not found.

(13) The properties of this soln are that it requires a large hash table. This is because hash table must be at least as large as set of keywords and designed to minimize the probability of collisions. This means the search operation can be performed quickly even for large set of keywords. This soln is also effective because it can find any keyword in set even if keyword is not stored in hash table.

The trade offs of this soln are that it requires a large hash tables.

This is because hash table must be at least as large as set of keywords...

However probability of this happening is very small as long as hash functions from universal hash functions family are used.

(iv) A universal hash function family was req. in solⁿ because it ensures that probability of collisions is minimized. This is req. because it ensures probability of collisions can significantly slow down the search operation. If a non-universal hash function used, the probability of collisions could be high especially for large set of keywords.

It is not possible to use any other hash function in this method without increasing probability of collisions. This is because the universal hash function family is designed to minimize the probability of collisions.

Any other hash functions would likely have a higher probability of collisions, which could make search operation slower.

(4)

A graph is connected if there is a path b/w any 2 vertices in graph. A connected graph does not consist of disconnected components; every vertex is reachable from every other vertex.

$G = (V, E)$ Here every pair of vertices u & v in V . There exists a path b/w u & v .

Here consecutive pair of vertices is connected with edge E .

necessary conditions:

- ② Graph must have atleast 2 vertices
- ① must be atleast one edge in graph

(ii) Process & steps

- 1) choose random subset of vertices in graph
- 2) Run algorithm on subset of vertices
- 3) if algo. o/p that the graph is connected, graph is considered connected else not connected.

(iii) Ex

$$\begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

choosing random subset of vertices $\{0, 1, 2\}$

- e if Algorithm graph is connected graphs considered as connected else not connected.
- if it is connected because of there is a path b/w subset of vertices .
- ∴ So Algo. correctly determines connectedness of graphs