

SS-ZG548: ADVANCED DATA MINING

08

Frequent Items Count Distinct



Dr. Kamlesh Tiwari

Assistant Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Sept 11, 2021

(WILP @ BITS-Pilani July-Dec 2021)

<http://ktiwari.in/adm>

Frequent items over data stream

- Let identity of items is drawn from the set $\{1, 2, 3, \dots, n\}$.
- Frequency of item i be f_i
- Assume general arrival model, (i, v) , $v > 0$ represents arrival and $v < 0$ is departure.
- Sum of frequencies $m = \sum_i f_i$ represent size of data stream
- Frequency** item i , have **frequency** $f_i > m/(k+1)$ for some fixed k

Observations

- There could be at most k frequent items (why ? proof? $m > k(k+1)$)
- Any algorithm that finds all frequent and only frequent items requires at least $\log \binom{n}{k}$ bits (how? $2^s \geq \binom{n}{k}$)

In action: Frequent items

Take $k = 4$, and consider following data stream

5	8	4	5	4	12	→
6	5	2	8	3	5	→
	4	5	4	12	6	13

Insert x in data structure

- IF (A.ismember(x)) A[x]++
- ELSE A.insert(x)
- IF (A.size == k+1) THEN $\forall y \in A$
A[y]--,
IF (A[y] ==0) A.delete(y);

Let us step by step execute the algorithm:

Index	Item	Frequency
1		
2		
3		
4		
5		

Recap: Data streams

Consider, stream of data. Where data in arriving in rapid succession. Re-scan is NOT possible. Even storage space is insufficient to accommodate all data points.

Without storing all the data one wish to estimate

- Set of frequent items
- Number of distinct items
- Frequent itemsets
- etc

from first n Natural numbers, without repetition, in an arbitrary order. Can you report the missing one?

[Constraints are on memory and processing power]

Find approx frequent items

Wish to output a list of items such that

- Every item in the list has frequency $f_i > (1 - \epsilon) \frac{m}{k+1}$
- All the items having frequency at least $(1 + \epsilon) \frac{m}{k+1}$ is in the list

Output should satisfy above two properties with probability $(1 - \delta)$

Algorithm maintains a data structure A over the stream.

Step to update an item x is as below

- IF (A.ismember(x)) A[x]++
- ELSE A.insert(x)
- IF (A.size == k+1) THEN $\forall y \in A$
A[y]--,
IF (A[y] ==0) A.delete(y);

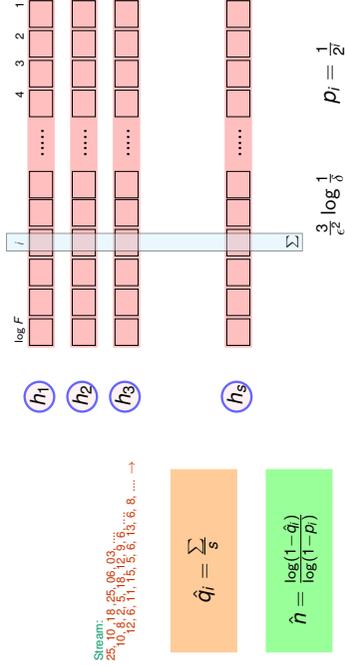
Count distinct over data streams (FM sketch)

Estimate number of distinct items in data stream

- If $x = ??? \dots ?1000\dots 0$ then $L[x]=i$
- Probability of $L[x]=i$ is $p_i = \frac{2^{\log_2 |F|-i}}{|F|} = 1/2^i$ when $x \in \{1, 2, \dots, F\}$
- FM sketch is a bitmap A of size $\log_2 |F|$ with hash a function h
- Arrival of an item x , sets bit $A[L[h(x)]] \leftarrow 1$.
Probability that $A[j] = 1$ after seeing n items is $1 - (1 - p_j)^n$
- With s independent copies of FM sketch, let $\#A[i]$ represent count of 1's at level i and $\hat{q}_i = \frac{\#A[i]}{s}$. Then choose i , such that $\hat{q}_i \geq \frac{3}{2} \log \frac{1}{\delta}$. By Chernoff's bound $\hat{n} \in [(1 - \epsilon)E[n], (1 + \epsilon)E[n]]$ with probability $(1 - \delta)$

$$\hat{n} = \frac{\log(1 - \delta)}{\log(1 - p_i)}$$

In action: Count distinct over data streams



$\hat{h} \in [(1 - \epsilon)E[n], (1 + \epsilon)E[n]]$ with probability $(1 - \delta)$

Thank You!

Frequent pattern mining over data streams

- Applications involves retail market data analysis, network monitoring, web usage mining, and stock market prediction.
- Using sliding window
- Efficiently remove the obsolete, old stream data
- Compact Pattern Stream tree (CPS-tree)
- Highly compact frequency-descending tree structure at runtime
- Efficient in terms of memory and time complexity
- Pane and window
- Insertion and restructuring

Thank you very much for your attention!

Queries ?