

13. Majority and frequent items

Friday, October 8, 2021 3:13 PM

Last time: Frequency moment sketches

Today: Majority & frequent item items

Majority item.

Ex n people voting for m candidates

Does any candidate have $> \frac{n}{2}$ votes?

i.e. Let $a_1, \dots, a_n \in [m]$. Determine if $\exists s \in [m]$ s.t. s occurs $> \frac{n}{2}$ times?

Claim: Any deterministic streaming alg requires $\Omega(\min(n, m))$ space, if we require that it output if there is a majority element, and if so, what!

proof. Suppose n is even and the last $\frac{n}{2}$ items are identical.

Every possible set of unique $\frac{n}{2}$ first items must have a different memory config, otherwise we can make a mistake by choosing second half to belong to one subset but not the other.

If $\frac{n}{2} \geq m$, then $2^m - 1$ subsets $\log(2^m - 1) = \Omega(m)$ bits

$\frac{n}{2} \leq m$, then $\geq \frac{m!}{(m - \frac{n}{2})!}$ subsets $\log\left(\frac{m!}{(m - \frac{n}{2})!}\right) = \Omega(n)$ bits



Majority Alg - with undefined behavior when no majority

Initialize $B \leftarrow a_1$ and $c \leftarrow 1$.

For a_i in $i \in \{2, \dots, n\}$,

If $B = a_i$, $c \leftarrow c + 1$

Else if $c > 0$, $c \leftarrow c - 1$

Else if $c = 0$, $B \leftarrow a_i$, $c \leftarrow 1$

} paired elim. of items

If $c > 0$, output B

If \exists majority item, appears $> \frac{n}{2}$ times, and so it cannot be eliminated.

Misra-Gries Algorithm Frequent

Initialize $B_1, \dots, B_k = 0$ buckets, and $c_1, \dots, c_k = 0$ counters

For $i \in [n]$

If $\exists j$ s.t. $B_j = a_i$, $c_j \leftarrow c_j + 1$

Else

If $\exists j$ s.t. $B_j = 0$, $B_j \leftarrow a_i$, $c_j \leftarrow c_j + 1$

Else (Decrement)

For all j

$c_j \leftarrow c_j - 1$

If $c_j = 0$, $B_j \leftarrow 0$.

Thm. 6.2 At the end of Misra-Gries, for each B_k with true count x_k ,

$$c_k \in \left[x_k - \frac{n}{k+1}, x_k \right]$$

If some $s \neq B_k$ for any k , then $x_k \leq \frac{n}{k+1}$.

Count-Min-Sketch

Recall Bloom filters. Probabilistic set membership query.

Maintain bit vector

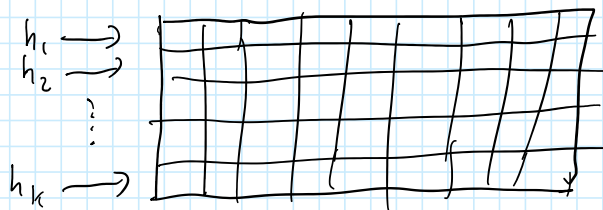


by setting bits corresponding to multiple ind. hash functions.

Query item s by checking if bits $h_1(s), \dots, h_k(s)$ are set.

Might accidentally return yes because of hash collisions, but unlikely given right params

Can also imagine each hash function having its own range.



Let's adapt this to freq. counting using Count-Min Sketch.

Let $a_1, a_2, \dots, a_n \in [m]$ and let $\vec{x} \in \mathbb{R}^m$ be the freq. of each item in $[m]$.

We want to estimate \vec{x} .

Maintain $t \times w$ matrix of counters

$$C = t \left\{ \begin{array}{c} \left[\begin{array}{c} h_1(j) \\ h_2(j) \\ \vdots \\ h_t(j) \end{array} \right] \\ w \end{array} \right.$$

For each row, associate a hash function $h_j: [m] \rightarrow [w]$ from a 2-wise family.

Insert item i by incrementing all counters $C_{j, h_j(i)}$ for $j \in [t]$.

Output PointQuery(i) = $\min_{j \in [t]} C_{j, h_j(i)}$

Claim: If $t > \log(\frac{1}{\delta})$ and $w \geq \frac{2}{\epsilon}$, then

$$\text{Prob} \left(\text{PointQuery}(i) \in [x_i - \epsilon \|\vec{x}\|_1, x_i + \epsilon \|\vec{x}\|_1] \right) \geq 1 - \delta$$

proof. For any $j \in [t]$,

$$C_{j, h_j(i)} = x_i + \sum_{\substack{r \neq i \\ h_j(r) = h_j(i)}} x_r = x_i + \sum_{r \neq i} \delta_r x_r$$

where $\delta_r = \begin{cases} 1 & \text{if } h_j(r) = h_j(i) \\ 0 & \text{otherwise} \end{cases}$

$$\mathbb{E} \sum_{r \neq i} \delta_r x_r = \frac{1}{w} \sum_{r \neq i} x_r \leq \frac{\epsilon}{2} \|\vec{x}\|_1$$

By Markov's inequality + since $x_i \geq 0$,

$$\text{Prob}(\text{noise} > \epsilon \|\vec{x}\|_1) \leq \frac{1}{2}.$$

So $C_{j, h_j(i)} \geq x_i$ and w.p. $> \frac{1}{2}$, $C_{j, h_j(i)} \leq \epsilon \|\vec{x}\|_1$,

So we are repeating that $t = \log(\frac{1}{\delta})$ times,

$$\text{Prob} \left(\min_{j \in [t]} C_{j, h_j(i)} > x_i + \epsilon \|\vec{x}\|_1 \right) = \text{Prob} \left(\forall j \in [t], C_{j, h_j(i)} > \epsilon \|\vec{x}\|_1 \right)$$

$$\text{Prob} \left(\min_{j \in [t]} \langle j, h_j(i) \rangle \cdot x_i + \sum \|x\|_1 \right) = \text{Prob} \left(\forall j \in [t], \langle j, h_j(i) \rangle \geq 2 \|x\|_1 \right) \leq \frac{1}{2^t} < \delta.$$

Only useful for heavy hitters (very frequent items) $x_i > \epsilon \|\vec{x}\|_1$,
 so useful for $\sim \frac{1}{\epsilon}$ of the values at most. □

Set-similarity

Let $A, B \subseteq U$ be two subsets. Let $n = |A \cup B|$.

Then Jaccard Similarity $J = J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ is a measure of set-similarity

MinHash (idealized)

1. Let $h_i: U \rightarrow [q]$ be ind. unif. random hash functions, $i \in [k]$

2. Let $\delta_i = \begin{cases} 1, & \min_{a \in A} h_i(a) = \min_{b \in B} h_i(b) \\ 0, & \text{else.} \end{cases}$

↗ using oracle actual random hash functions

3. Then output $\hat{J} = \frac{1}{k} \sum_{i=1}^k \delta_i$ as estimate for J .

Claim: If $q > \frac{kn^2}{\delta}$, and $k > \frac{2}{\epsilon^2 \delta}$, then

$$\text{Prob} \left(|\hat{J} - J| > \epsilon \right) < \delta.$$

proof. Recall that if h_i is a universal hash function, all of $h_i(x), x \in A \cup B$ are distinct with probability at least $1 - \frac{\delta}{2k}$. (signatures)

By union bound, all h_i have no collisions w.p. at least $1 - \frac{\delta}{2}$.

Then $h_i(a) = h_i(b)$ only if $a = b$.

$$\Rightarrow \min_{a \in A} h_i(a) = \min_{b \in B} h_i(b) \text{ only if } a = b \in A \cap B.$$

Clearly, the converse also holds, so $\mathbb{E} \delta_i = J$, and $\mathbb{E} \hat{J} = J$, $\text{Var}(\delta_i) \leq J$.

$$\text{Then } \text{Prob} \left(|\hat{J} - J| \geq \epsilon \right) \leq \frac{\text{Var}(\hat{J})}{\epsilon^2} \leq \frac{1}{\epsilon^2} \cdot \frac{1}{k} \cdot J < \frac{\delta J}{2} \leq \frac{\delta}{2}. \quad \square$$

Thus, we need $O\left(\frac{1}{\epsilon^2} \log \frac{kn^2}{\delta}\right) = O\left(\frac{1}{\epsilon^2} \log \frac{n}{\epsilon \delta}\right)$
not tight at all.

As.2e: Can do better in idealized setting $O\left(\frac{1}{\epsilon^2} \left(\log \log n + \log \frac{1}{\epsilon}\right)\right)$
[Yu, Weber, TKDE, 2020] for ϵ additive error

As.2c: We assume hash functions were fully random

Unfortunately, cannot just use 2- or 4-wise hash families,
but need a stronger cond., i.e., that any item is equally
likely to be the minimum.

A new criterion called min-wise independence [Indyk, 1999]
↳ almost achievable for l -wise ind.
hash functions for l large enough.