

CS7960 L11 : Streaming | Approximate Rank

Streaming Algorithms

Stream : $A = \langle a_1, a_2, \dots, a_m \rangle$
 $a_i \in [n]$ size $\log n$
Compute $f(A)$ in $\text{poly}(\log m, \log n)$ space

Let $f_j = |\{a_i \in A \mid a_i = j\}|$

$F_1 = \sum_j f_j = m == \text{total count}$

Goal:

ϵ -RELATIVE-RANK: Build data structure S .
 $\text{rank}(v) = 1 + \# \text{ items in } A \text{ smaller than } v$
 $\text{relative-rank}(v) = Rrank(v) = \text{rank}(v)/|A| \text{ in } [0,1]$

ϵ -RELATIVE-RANK S returns $S(v)$ such that
 $Rrank(v) - \epsilon \leq S(v) \leq Rrank(v) + \epsilon$

Why relative rank?

$\text{median}(A) = \{ * \mid Rrank(*) = 1/2 \}$
quantiles: $Rrank(*) = 1/4, 3/4$
- more "robust" than mean
- captures distribution

Warm Up:

Random Sampling | Reservoir Sampling S w/ $k = O(1/\epsilon^2 \log(1/\delta))$
For any interval $I_v = (-\infty, v)$ --> $Rrank(v) = \text{count}(I)/m$
Estimate $\text{count}(I_v)$ w/ $S(v)$ s.t.
 $Rrank(v) - \epsilon < S(v)/m < Rrank(v) + \epsilon$
correct w/ $1-\delta$

Want S of size $\sim O(1/\epsilon)$ instead of $\sim O(1/\epsilon^2)$

ϵ		.1	.01	.001
$1/\epsilon$		10	100	1000
$1/\epsilon^2$		100	10000	1000000

Greenwald-Kanna Algorithm [G,K '01]

```

min-rank(v) = smallest possible rank of v (to S's knowledge)
max-rank(v) = largest possible rank of v (to S's knowledge)
m = size of stream up to know
| A = <a1, a2, ..., a_m | a_{m+1}, ... >

Maintain collection of tuples (v_i, g_i, Delta_i)
s.t. v_i <= v_{i+1}
     g_i = min-rank(v_i) - min-rank(v_{i-1})
     Delta_i = max-rank(v_i) - min-rank(v_i)
Maintain MIN = min{a_j \in A} (seen so far)
     MAX = max{a_j \in A} (seem so far)
Maintain m = total count

#####
Process: a_m
Find i s.t. v_i < a_m < v_{i+1}
a_m becomes v_{i+1} (don't actually keep indices)
(v_{i+1} = a_m, g_{i+1} = 1, Delta_i = L 2 eps m J )
m <- m+1
Update {MIN, MAX} as needed
--> Delta_i = 0
Check if we can compress:
if (i s.t. g_i + g_{i+1} + Delta_{i+1} <= L 2 eps m - 1J) then
    Remove (v_i, g_i, Delta_i)
    Set g_{i+1} <- g_i + g_{i+1}
#####

Query : rank(v)?
Find i s.t. v_i < v < v_{i+1}
Return (min-rank(v_i) + max-rank(v_{i+1}))/2
- min-rank(v_i) = 1 + sum_{j=1}^i g_i
- max-rank(v_{i+1}) = min-rank(v_{i+1}) + Delta_{i+1}

Error ?
- we know : max-rank(v_{i+1}) - min-rank(v_i) = g_i + g_{i+1} + Delta_{i+1}
g_i + g_{i+1} + Delta_{i+1} <= 1 + L 2 eps m - 1 J < 2 eps m
INDUCTION: Base case true on insertion
            Delta < 2 eps m, by induction as well...
            Insertion does not increase
            Compression only allowed if above holds
            2 eps m / 2 < eps m -> Rrank(v) within eps.


```

Space : $O((1/\epsilon) \log(\epsilon * m) \log n)$

- each tuple space : $O(\log n + \log m) \rightarrow O(\log n + \log (\epsilon * m))$
- how many tuples? $O((1/\epsilon) \log (\epsilon * m))$

INTUITION: <analysis quite complicated>

- Delta starts 2 eps m , but as m grows, this shrinks relative to m .
- g_i starts at 1. If $g_i + g_{i+1} < \text{eps } m$, eventually will compress as Δ_{i+1} decreases relative to m .
- So $g_i + g_{i+1} > \text{eps } m$ unless g_{i+1} is new
(second half of stream).

<this is hard part to show,
need amortized slack of $\log(\text{eps } m)$ >

- if $g_i + g_{i+1} > \text{eps } m$, then $g_i > i * \text{eps } m / 2$
 $\rightarrow \max i < 2 / \text{eps}$.

So $O((1/\text{eps}) * \log(\text{eps } m)) * O(\log n + \log(\text{eps } m))$
<second $\log(\text{eps } m)$ can be absorbed into
the $O((1/\text{eps}) \log(\text{eps } m))$ term>
--> $O((1/\text{eps}) \log(\text{eps } m) \log n)$

Median Exactly?

1-pass $\Omega(\min\{m, n\})$ space
2-pass $\tilde{O}(\sqrt{n})$ space
 p -pass $\tilde{O}(n^{1/p})$ space

$p = O(\log n)$ passes --> $\tilde{O}(1)$ space

2-pass set $\text{eps} \sim 1/\sqrt{n}$ in GK.

- Space: $O(\sqrt{n} * \log(m/\sqrt{n}) \log n)$
- Estimate rank of all elements within \sqrt{n} on pass 1
Find $\max i$ s.t. $R\text{rank}(v_i) < 1/2$
 $\max j$ s.t. $R\text{rank}(v_j) > 1/2$
- $I_{i,j} = (v_i, v_j)$ and $\text{count}(I_{i,j}) = O(\sqrt{n})$
- Second pass keep all elements in $I_{i,j}$ and count a in A s.t. $a < v_i$

In p -pass narrow range and approximate each pass...

A few more passes can get much more accurate!