

# Lec31-count-distinct-and-MinHash

Sunday, December 3, 2023 7:35 PM

## Distinct - element sketching

Consider sequence  $a_1, a_2, \dots, a_n \in [m]$ , with  $n, m \gg 0$ .

Want to know number of unique items,  $F_0 = |\text{uniq}\{a_1, \dots, a_n\}|$ ,

$O(m)$  space solution - bit vector

$O(n \log m)$  space sol - store list of items seen

Goal:  $O(\log n \cdot \log m)$  space.

Thm: (lower bound). Any exact deterministic algorithm must use at least  $m$  bits of memory on some seq. of length  $m+1$ .

proof: Assume that an alg uses  $< m$  - bits of memory on all such seq.

Recall that the power set  $|\mathcal{P}([m])| = 2^m$  and  $\text{uniq}\{a_1, \dots, a_m\}$  can be any subset except  $\emptyset$ , so  $2^m - 1$  possibilities.

But, we only have  $2^{m-1}$  memory states possible since we use  $< m$  bits.

Thus, two diff. subsets  $S_1, S_2 \in \mathcal{P}([m]) \setminus \emptyset$  must have same memory state.

$\Rightarrow |S_1| = |S_2|$  because otherwise our alg would be wrong.

Let  $b \in S_1$  but  $b \notin S_2$ . Then  $S_1 \cup \{b\} = S_1$  so  $|S_1 \cup \{b\}| = |S_1|$   
and  $S_2 \cup \{b\} \neq S_2$ , so  $|S_2 \cup \{b\}| = |S_2| + 1$ .

But alg will have same mem state after adding  $b$ , so  
it must be wrong on one of them



## Consider Idealized Streaming Algorithm (ISA)

1. Pick random hash function  $h: [m] \rightarrow [0, 1]$

(where  $h(x) \sim \text{Unif}[0, 1]$  iid)

2. Calculate  $z = \min_{i \in \text{stream}} h(a_i)$

3. Output  $\frac{1}{z} - 1$ .

Let  $\mathcal{S} = \text{unig} \{a_1, \dots, a_n\} = \{b_1, \dots, b_d\}$   $\leftarrow$   $d$  distinct elements

$h(b_1), \dots, h(b_d) = X_1, \dots, X_d$  i.i.d.  $\text{Unif}[0, 1]$

$$Z = \min_{i=1}^d \{X_i\}$$

Lemma: If  $X: \Omega \rightarrow [0, +\infty)$  is a nonnegative r.v., then

$$\mathbb{E}X = \int_{[0, +\infty)} \text{Prob}(X > x) dx$$

Claim:  $\mathbb{E}Z = \frac{1}{d+1}$

proof:  $\mathbb{E}Z = \int_0^\infty \text{Prob}(Z > \lambda) d\lambda = \int_0^1 \text{Prob}(\forall i, X_i > \lambda) d\lambda$

$$= \int_0^1 \prod_{i=1}^d \text{Prob}(X_i > \lambda) d\lambda = \int_0^1 (1-\lambda)^d d\lambda = \left[ \frac{-(1-\lambda)^{d+1}}{d+1} \right]_0^1 = \frac{1}{d+1}$$

Claim:  $\mathbb{E}Z^2 = \frac{2}{(d+1)(d+2)}$

proof:  $\mathbb{E}Z^2 = \int_0^1 \text{Prob}(Z^2 > \lambda) d\lambda = \int_0^1 \text{Prob}(Z > \sqrt{\lambda}) d\lambda$

$$= \int_0^1 (1-\sqrt{\lambda})^d d\lambda = 2 \int_0^1 u^d (u-1) du = 2 \int_0^1 (u^{d+1} - u^d) du$$

let  $u = 1 - \sqrt{\lambda}$   
 $u^2 - 2u + 1 = \lambda$   
 $(2u-2)du = d\lambda$

$$= 2 \left[ \frac{1}{d+2} u^{d+2} - \frac{1}{d+1} u^{d+1} \right]_0^1$$

$$= 2 \left[ \frac{-1}{d+2} + \frac{1}{d+1} \right] = \frac{2}{(d+1)(d+2)}$$

Thus,  $\text{Var}(Z) = \mathbb{E}Z^2 - (\mathbb{E}Z)^2 = \frac{d}{(d+1)^2(d+2)} < \frac{1}{(d+1)^2}$ .

Averaging algorithm (AA)

1. Run  $a = \frac{1}{n}$  ISAs in parallel

1. Run  $q = \frac{1}{\epsilon^2 \eta}$  ISAs in parallel.

2.  $\bar{z} = \frac{1}{q} \sum_{i=1}^q z_i$ , where  $z_i$  comes from ISA.

3. Output  $\frac{1}{\bar{z}} - 1$ .

$$\text{Then } E(\bar{z}) = \frac{1}{d+1}, \quad \text{Var}(\bar{z}) = \frac{1}{q} \cdot \frac{d}{(d+1)^2(d+2)} < \frac{1}{q(d+1)^2}$$

$$\text{By Chebyshev, } \text{Prob}\left(\left|\bar{z} - \frac{1}{d+1}\right| > \frac{\epsilon}{d+1}\right) < \frac{(d+1)^2}{\epsilon^2} \cdot \frac{1}{q(d+1)^2} = \eta.$$

Claim:  $\text{Prob}\left(\left|\frac{1}{\bar{z}} - d - 1\right| > O(\epsilon)d\right) < O(\eta).$

proof:  $\text{Prob}\left(\left|\bar{z} - \frac{1}{d+1}\right| > \frac{\epsilon}{d+1}\right) < \eta \iff \text{w.p. } 1 - \eta, \left|\bar{z}\right| \leq \frac{1 + \epsilon}{d+1}$

$$\text{Prob}\left(\left|\bar{z}d + \bar{z} - 1\right| > \epsilon\right) < \eta$$

$$\text{Prob}\left(\left|\frac{1}{\bar{z}} - d - 1\right| > \frac{\epsilon}{|\bar{z}|}\right) < \eta$$

$$\Rightarrow \text{Prob}\left(\left|\frac{1}{\bar{z}} - d - 1\right| > \underbrace{\frac{\epsilon(d+1)}{1+\epsilon}}_{O(\epsilon)d} \underbrace{\frac{1}{1+\epsilon}}_{O(\eta)}\right) < \eta$$



With high prob, our estimator is within a factor of  $1 + O(\epsilon)$  of  $1/d$ .

Space complexity  $O\left(\frac{1}{\epsilon^2 \eta} \cdot \text{space}[0, 1]\right)$ .

Can do even better analysis, but this is core of idea of  
loglog family of estimators.

Can do even better analysis, but this is cool

LogLog family of estimators.

### 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  $h_i: U \rightarrow [k]$  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  $\epsilon > \frac{k n^2}{8}$ , and  $k > \frac{2}{\epsilon^2 \delta}$ , then

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

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  $\mathbb{E} \hat{J} = J$ ,  $\text{Var}(\delta_i) \leq J$ .  
Then  $\text{Prob}(|\hat{J} - J| \geq \epsilon) \leq \frac{\text{Var}(\hat{J})}{\epsilon^2} \leq \frac{1}{\epsilon^2} \cdot \frac{1}{k} \cdot J \leq \frac{\delta J}{2} \leq \frac{\delta}{2}$ . □

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

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

Aside: We assume hash functions were fully random

Unfortunately, cannot just use 2-wise 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, 1997]

↳ almost achievable for  $l$ -wise ind.  
hash functions for  $l$  large enough.