

L14: Distinct Item Counting

+ Mergable Summaries

Oct 13, 2025
Data Mining
Jeff M. Phillips



Streaming Input $A = \langle a_1, a_2, \dots, a_m \rangle$

observe in order, one pass

small space

$\mathcal{O}(\text{poly}(\log m, \log n))$

$a_i \in [n]$

count

label.

$$f_j = |\{a_i \in A \mid a_i = j\}| \quad \text{frequency}$$

$$F_0 = \# \text{ distinct items } j \in [n]$$

$\# f_j > 0$

↳ need compact way to maintain unique representations

Leverage hash functions

$$h : [n] \rightarrow [k]$$

↑ space of A

Bloom Filters

$$t \approx \frac{k}{|S|}$$

Maintain which items \leq in A

- no false negatives
- allow false positives

t hash functions $h_1, h_2, \dots, h_t \sim_{\text{iid}} H$

have array $B[0, \dots, k-1]$ of bits, init 0.

For $a_i \in A$
for $j = 1$ to t
Set $B[h_j(a_i)] = 1$

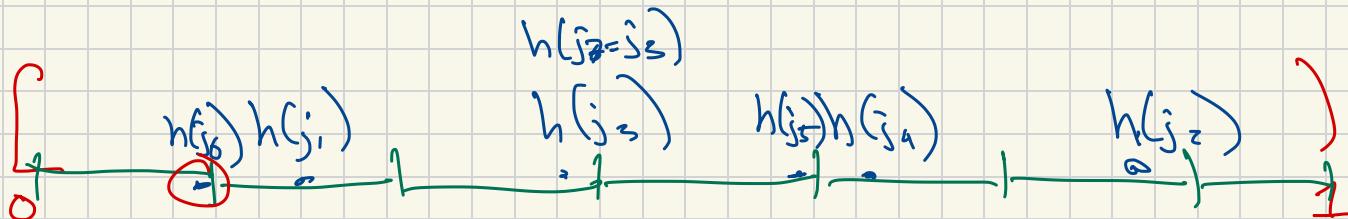
| $\begin{array}{l} \text{is } g \in S \\ \text{check if } h_j(g) = 1 \\ \forall j \in 1 \dots t \end{array}$

hash function

$$h: [n] \rightarrow [0, 1) \subset \mathbb{R}$$

continuous

how do we use to rep a round?



- each $h(j)$ random $\text{Unif}(0, 1)$

- number instances = # distinct from F_0

$$\gamma = \min_{a_i \in A} h(a_i)$$

$$\hat{F}_0 = \frac{1}{\gamma} - 1$$

$$E[\hat{F}] = F_0$$

Amplification

of hash functions

$$h_1, h_2, \dots, h_t \text{ iid}$$

$$\gamma = \frac{1}{t} \sum_i \delta_{h_i}$$

$$\bar{F}_0 = \frac{1}{\gamma} - 1$$

$$t = \frac{1}{\epsilon^2} \frac{1}{\delta}$$

w.p. $(1 - \delta)$

$$\left| \bar{F}_0 - F_0 \right| \leq \epsilon F_0$$

$$\epsilon = \frac{1}{10}$$

$$\delta = \frac{1}{10}$$

Issue #1: Failure prob.
high.

$$\bar{F} = F_0 + 10\% F_0$$

Issue #2: continuous, very small value γ .

Amplifies prob failure δ .

Median Trick

Run before

total size

$$q = \xi$$

$$\mathcal{O}(\gamma_q^2)$$

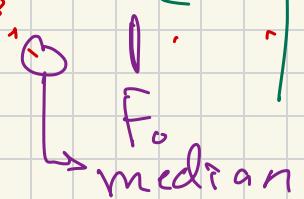
$$\delta = \gamma_2$$

no dependence of δ .

$$T = 2 \log \frac{1}{\delta}$$

Copies

$$\xi F_0 \quad \xi F_0$$



median T instances

Hyper Log Log

Flajolet - Martin

From continuous

$\gamma \rightarrow$ bits.

$$h(a_i) \rightarrow 0101101 \in \{0, 1\}$$

$$h(a_2) \rightarrow 00001000 = \frac{1}{16} + \frac{1}{128}$$

$\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}, \frac{1}{64}, \frac{1}{128}$

$$h(a_3) \rightarrow 01101111$$

$$h(a_4) \rightarrow 10-----$$

need $\log_2(F_0)$ bits.

just first bit $\Rightarrow \approx$ approx

storing the first 1 bit, requires $\log \log(F_0)$ bits

Mergeable Summaries

Summary as an S

- insert 1 item
- return (approx) query

$$\rightarrow \bullet \text{ merge } (S_A, S_B) \Rightarrow S_{A \cup B}$$

A 1 stream

B 1 disjoint streams

space $S(\frac{1}{\epsilon})$

$$\text{e.g. } S(\frac{1}{\epsilon}) = \frac{1}{\epsilon} \text{ MG, CS}$$

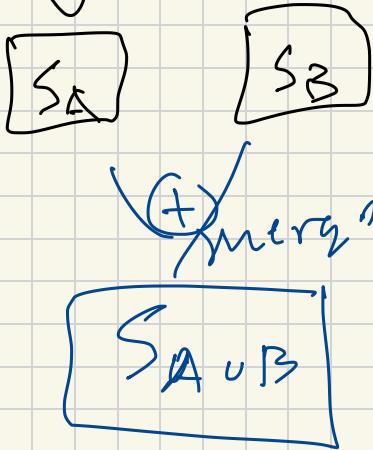
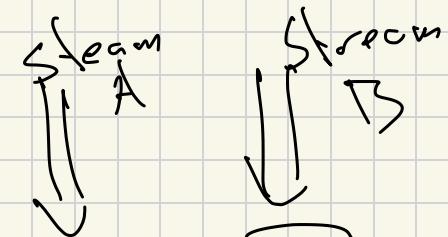
$$S(\frac{1}{\epsilon}) = \frac{1}{\epsilon^2} \text{ sample, F_0}$$

want $S_{A \cup B}$ same $S(\frac{1}{\epsilon})$?

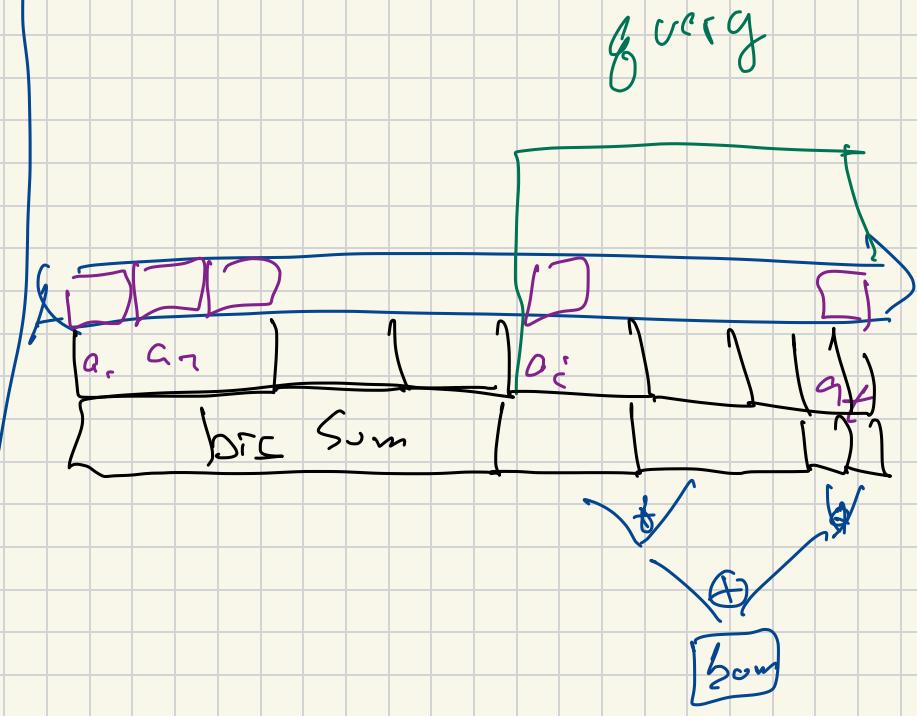
Approx Data Structures

→ no increase
in space
as error ϵ .

with Margins



in parallel



Mergeable

Sampling

each $a_i \rightarrow v_i \sim \text{Unif}(0, 1)$

Maintain $\min v_i \Rightarrow q_i$:

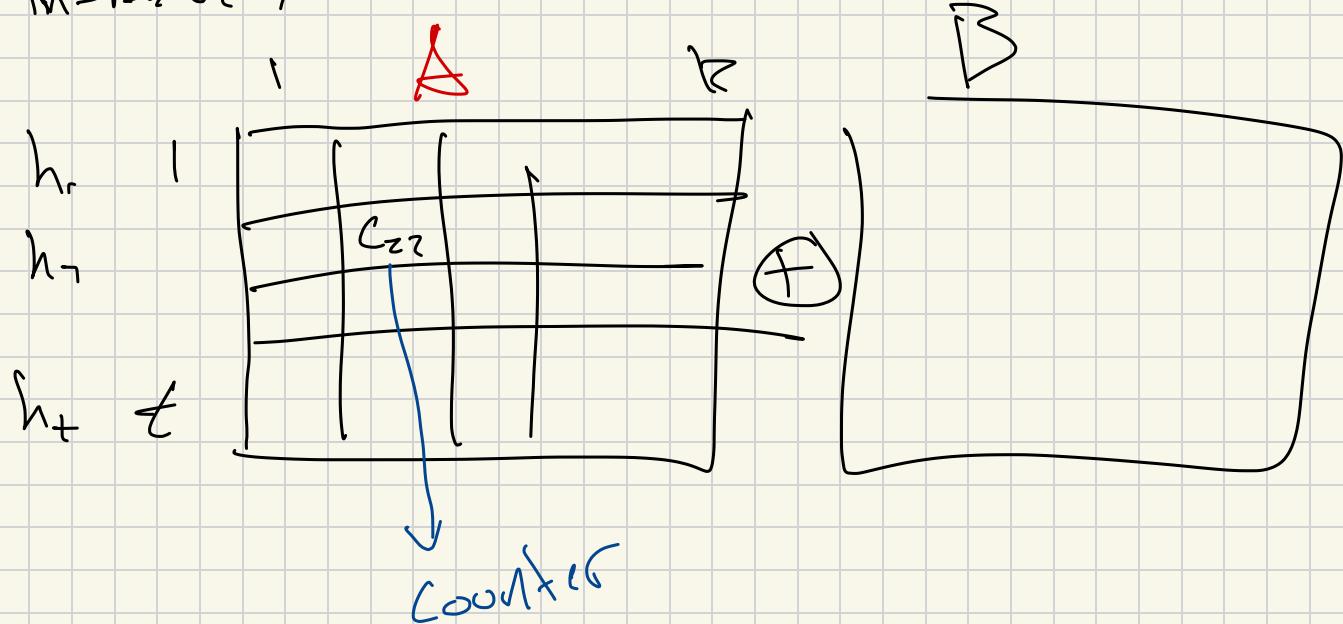
$$A = \begin{bmatrix} v_1 & v_2 & v_3 \\ a_1 & a_2 & a_3 \end{bmatrix}$$

$$B = \begin{bmatrix} v_1' & v_2' & v_3' \\ b_1 & b_2 & b_3 \end{bmatrix} \Rightarrow A \cup B$$

3 smallest

$$\begin{bmatrix} v_1' & v_1 & b_2' \\ \downarrow. & a_1 & b_2 \end{bmatrix}$$

Count Mismatch



add all counters

f_0 - stat. inh

Flugolist

Hyperr Log Log

$$t \cdot T = \frac{1}{\varepsilon^2} \log \frac{1}{\delta} \quad \text{of} \quad \begin{cases} \max_{i=1}^n \text{first}(k_i) \\ \dots \end{cases}$$

merges \rightarrow \hookrightarrow free max.