

CountMin Sketch: Finding Heavy Hitters

COMPCSI 753: Algorithms for Massive Data

Instructor: Ninh Pham

University of Auckland

Auckland, Aug 24, 2020

Basic definitions

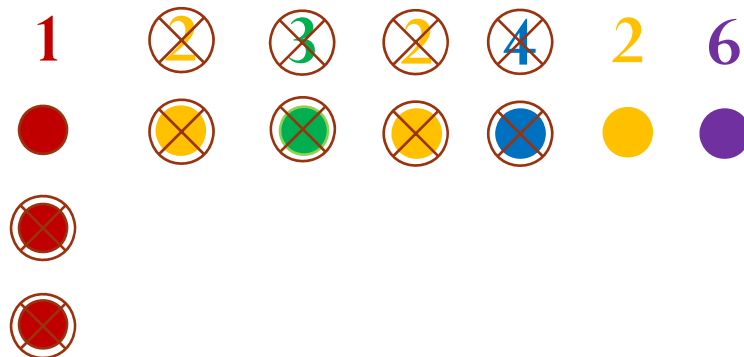
- Let \mathbf{U} be a universe of size \mathbf{n} , i.e. $\mathbf{U} = \{1, 2, 3, \dots, \mathbf{n}\}$.
- Cash register model stream:
 - Sequence of \mathbf{m} elements $\mathbf{a}_1, \dots, \mathbf{a}_m$ where $\mathbf{a}_i \in \mathbf{U}$.
 - Elements of \mathbf{U} may or may not occur once or several times in the stream.
- Finding heavy hitters in data stream (today's lecture):
 - Given a stream, finding frequent items.

Frequent items

- Each element of data stream is a tuple.
- Given a stream of m elements $\mathbf{a}_1, \dots, \mathbf{a}_m$ where $\mathbf{a}_i \in \mathbf{U}$, finding the most/top- k frequent elements.
- Example:
 - $\{\underline{1}, 2, \underline{1}, 3, 4, 5\} \rightarrow \mathbf{f} = \{\underline{2}, 1, 1, 1, 1\}$
 - $\{\underline{1}, \underline{2}, \underline{1}, 3, \underline{1}, \underline{2}, 4, 5, \underline{2}, 3\} \rightarrow \mathbf{f} = \{\underline{3}, \underline{3}, 2, 1, 1\}$
- We need an approximation solution with much smaller memory with theoretical guarantees.

Deterministic: Misra Gries

- Process an element **a**:
 - If we already have a counter for **a**, increment it.
 - Else, if there is no counter for **a**, but fewer **k** counters, create a counter for **a** initialized to 1.
 - Else, decrease all counters by 1. Remove 0 counters (key step).
- Example: $\{1, 2, 3, 1, 4, 2, 1, 4, 5, 2, 6\}$, $n=6$, $k=3$, $m=11$
 $\{1, 2, 3, 1, 4, 2, 1, 4, 5, 2, 6\}$



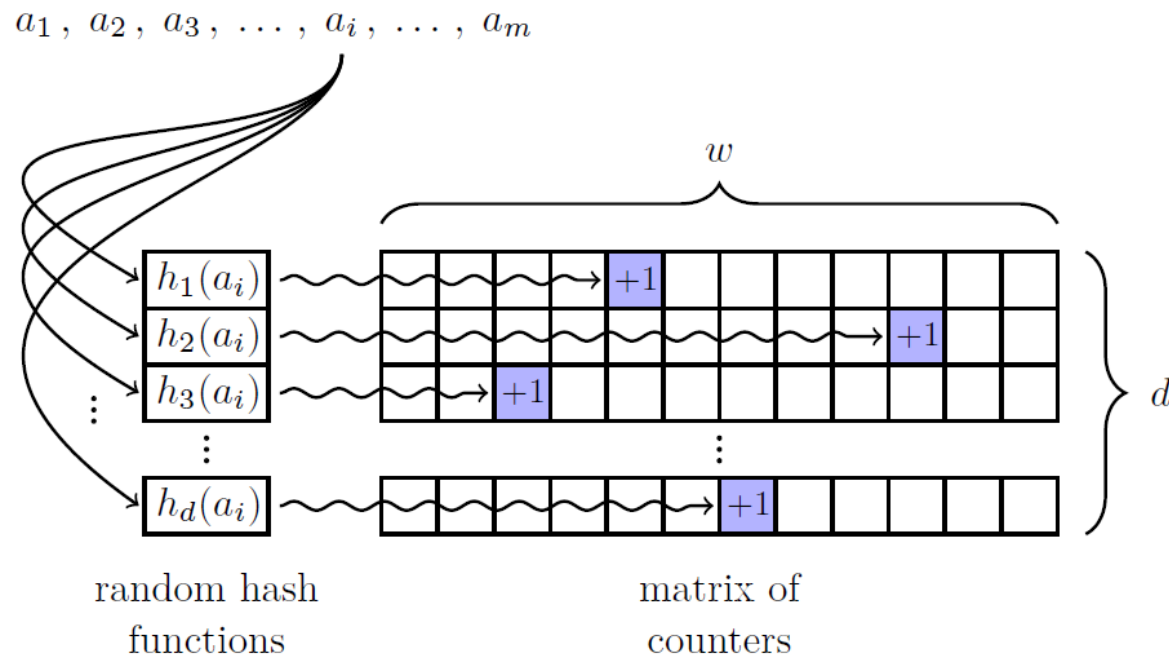
Randomized: CountMin sketch

- Setup:

- d independent **universal** hash functions \mathbf{h} over range $[0, w)$
- d different counters, C_1, \dots, C_d . Each of size w initialized with 0s.

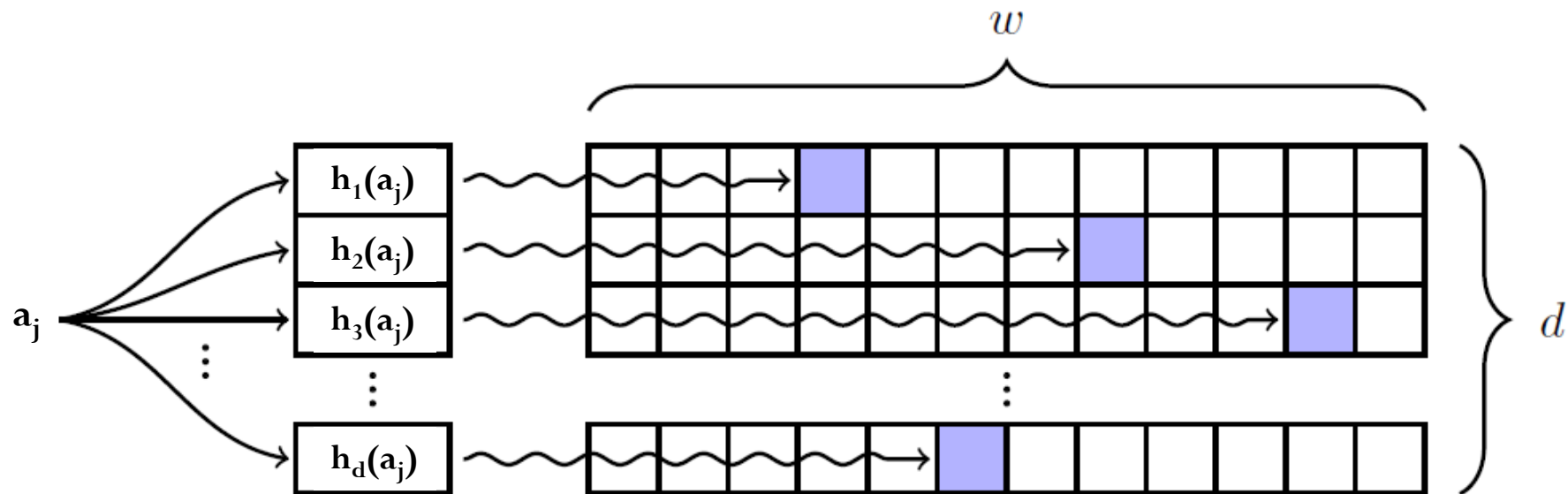
- Process an element \mathbf{a}_j :

- For each hash function, compute $\mathbf{h}_i(\mathbf{a}_j)$ and increment $C_i[\mathbf{h}_i(\mathbf{a}_j)]$ by 1



Randomized: CountMin sketch

- Query: How many times a_j occurred?
 - For each hash function, compute $h_i(a_j)$ and get $C_i[h_i(a_j)]$
 - Return $\min(C_1[h_1(a_j)], \dots, C_d[h_d(a_j)])$



return the minimum of values in blue cells

Universal hash function family

- Universal hash function definition:

- A family of hash function $H = \{h : U \rightarrow \{0, 1, \dots, w-1\}\}$ is **universal** if for any 2 distinct keys $x_i \neq x_j \in U$, we have

$$\Pr_h[h(x_i) = h(x_j)] \leq 1/w$$

- In our CountMin sketch:

- Given two different items $a_i \neq a_j$, what is the prob. a_i and a_j collide?

$$\Pr_h[h(a_i) = h(a_j)] \leq 1/w$$

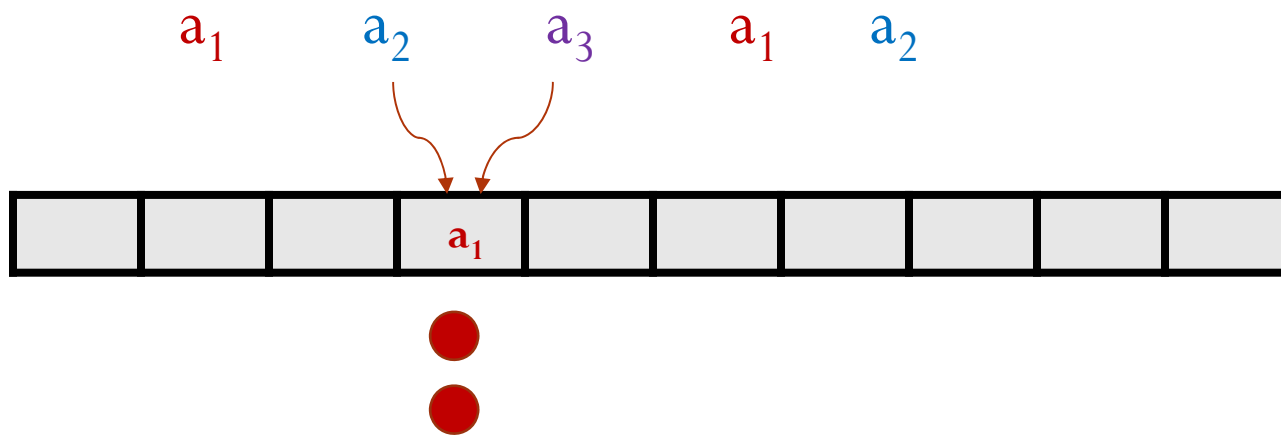
Analysis on 1 array of counters

- Notation:

- Stream of m items $\{a_1, \dots, a_m\}$ from the universe \mathbf{U} of size n .
- Frequency vector $\mathbf{f} = \{f_1, \dots, f_n\}$ and $\|\mathbf{f}\|_1 = m$.

- Question:

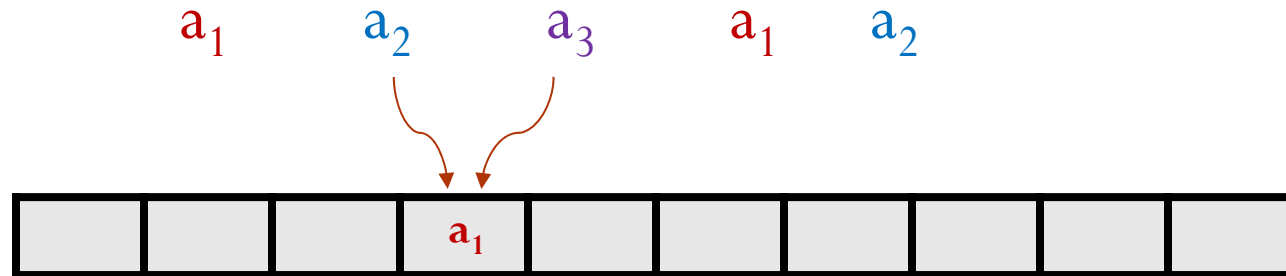
- Given a particular item a_1 , how many times $a_i \neq a_1$ collide by h ?



Analysis on 1 array of counters

- Question:

- Given a particular item a_1 , how many times $a_i \neq a_1$ collide by h .



- Answer:

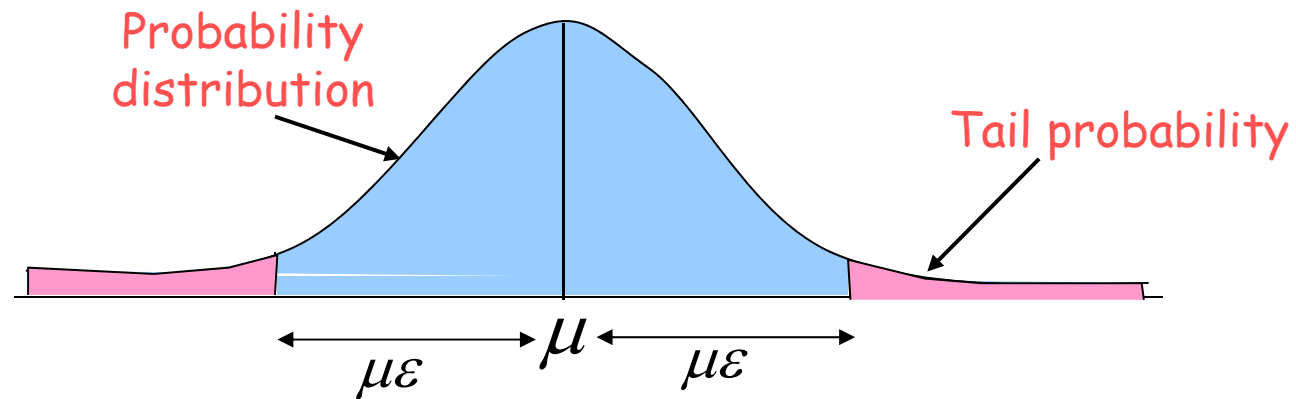
- Let X_2, \dots, X_n be contributions of a_2, \dots, a_n in the bucket $h(a_1)$.

$$\Pr [h(a_i) = h(a_1)] \leq 1/w \rightarrow E [X_i] \leq f_i/w$$

- Let $Y = X_2 + \dots + X_n$ be the total increments by other item $a_i \neq a_1$.

$$E[Y] = E[X_2] + \dots + E[X_n] = (f_2 + \dots + f_n)/w \leq m/w$$

Basic tools: Tail inequality



- Markov's inequality for $\mathbf{E}[Y] = \mu$:

$$\Pr[Y \geq 1 + \epsilon] \leq \frac{\mu}{1 + \epsilon} \text{ for any } \epsilon > 0.$$

$$\Pr[Y \geq (1 + \epsilon) \mu] \leq \frac{1}{1 + \epsilon} \text{ for any } \epsilon > 0.$$

Analysis on 1 array of counters

- **Observation:** We always over-estimate.
- **Question:** How large we over estimate?
- **Analysis for a particular item a_1 :**
 - Let $Y = X_2 + \dots + X_n$ be the total increments by other item $a_i \neq a_1$.
 - The value of counter $h(a_1)$: $f'_1 = f_1 + Y$.
 - Using Markov's inequality:
$$\Pr [f'_1 \geq f_1 + \epsilon m] = \Pr [Y \geq \epsilon m] \leq E[Y] / \epsilon m \leq m / (w * \epsilon m)$$
 - Choose $w = 2/\epsilon$, we have this prob. is at most $1/2$.
- Choose $w = 2/\epsilon$, the probability that our error is larger than ϵm is smaller than $1/2$.

Analysis on d arrays of counters

- Boosting the accuracy:

- Using d independent hash functions corresponding to d independent arrays of counters.
- $F_1 = \min(C_1[h_1(a_1)], \dots, C_d[h_d(a_1)]) = \min(f'_1, f'_2, \dots, f'_d)$.

- Analysis:

- $\Pr [F_1 \geq f_1 + \epsilon m] = \Pr [f'_1 \geq f_1 + \epsilon m \wedge \dots \wedge f'_d \geq f_1 + \epsilon m] \leq 1/2^d$
- Choose $d = \log(1/\delta)$, we have $\Pr [F_1 \geq f_1 + \epsilon m] \leq \delta$.

- With probability at least $1 - \delta$, we have

$$F_1 < f_1 + \epsilon m = f_1 + \epsilon \|f\|_1$$

Homework

- Implement the CountMin Sketch algorithm on the dataset from assignment 1:
 - **Description:** Each line (doc ID, word ID, freq.) as a stream tuple.
 - **Query:** What are the most and top-**10** frequent word ID have been used?