

Count 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.

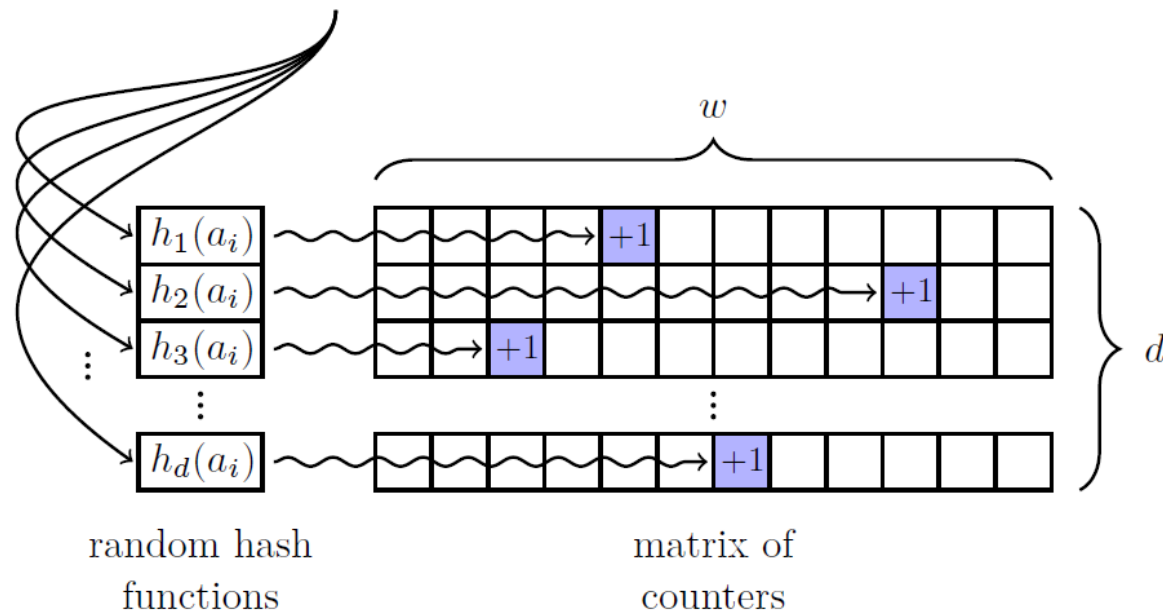
Randomized: CountMin sketch

- Setup:

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

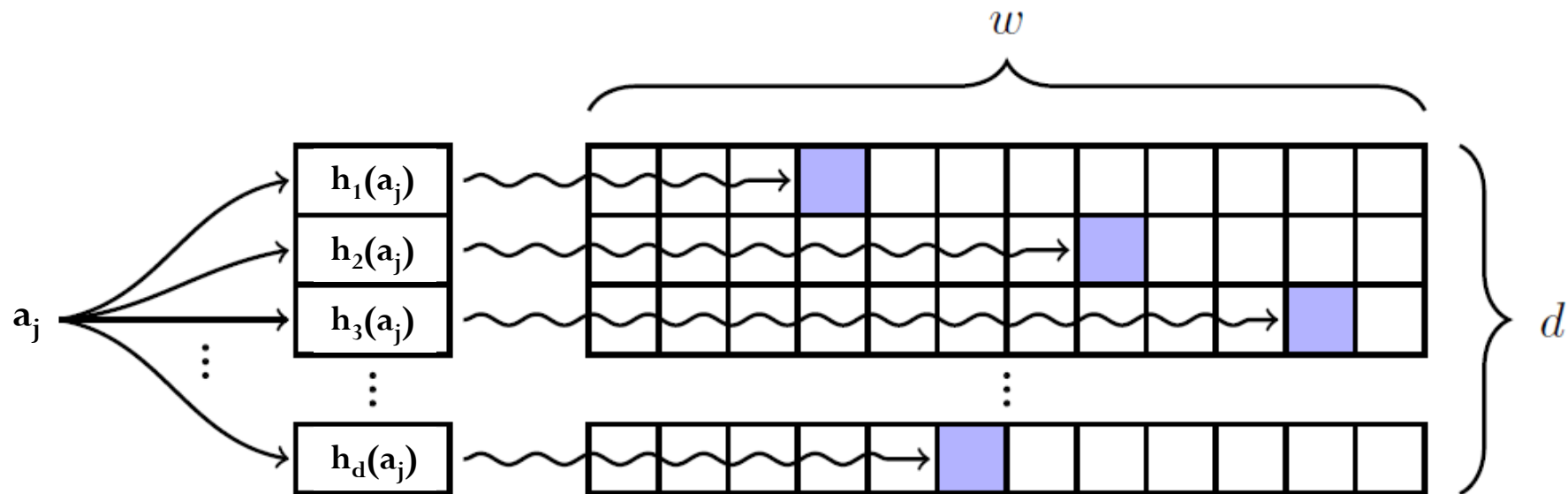
- Process an element a_j :

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

$$a_1, a_2, a_3, \dots, a_i, \dots, a_m$$


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

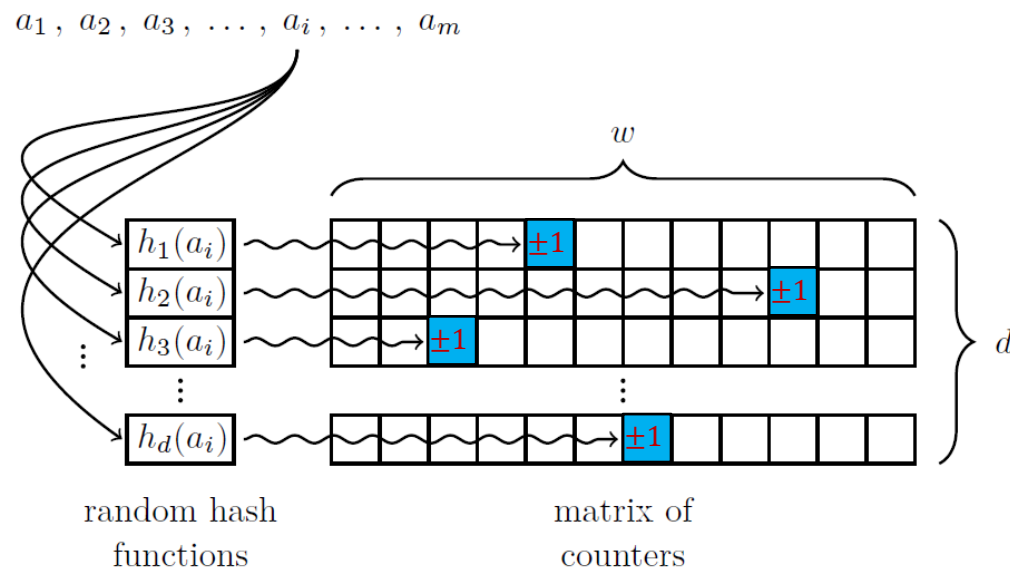
Randomized: Count sketch

- Setup:

- d independent **2-wise** hash functions \mathbf{h} over range $[0, w)$.
- d independent **2-wise** hash functions \mathbf{s} over range $\{+1, -1\}$.
- d different counters, $\mathbf{C}_1, \dots, \mathbf{C}_d$. Each of size w initialized with 0s.

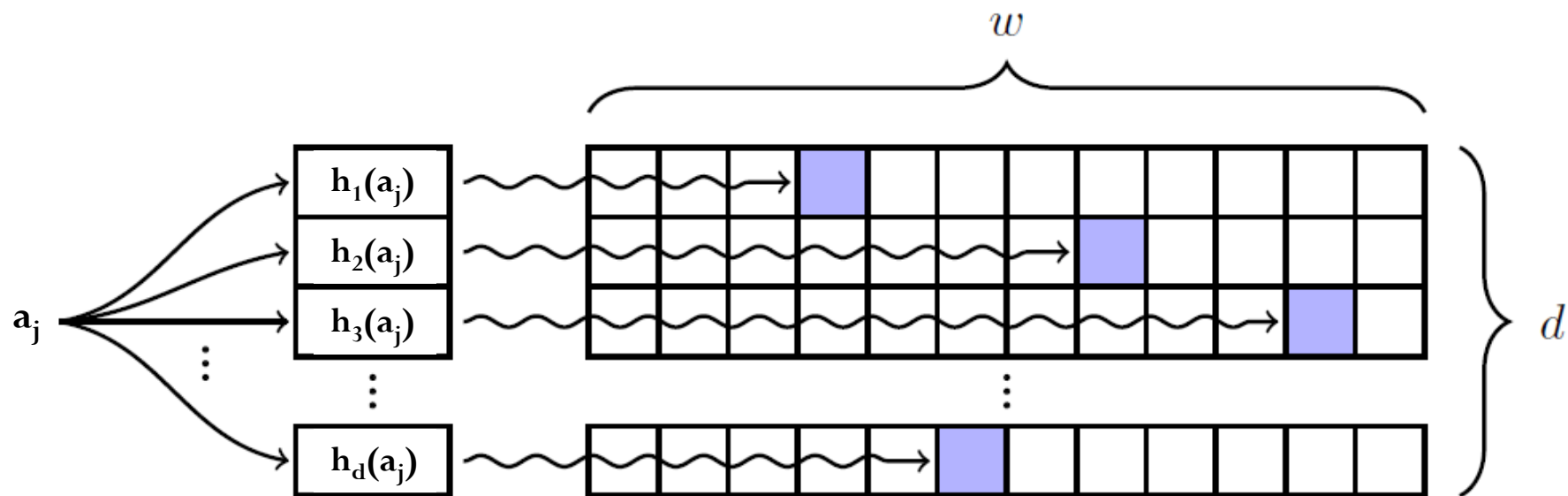
- Process an element \mathbf{a}_j :

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



Randomized: Count 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 **median**($C_1[h_1(a_j)], \dots, C_d[h_d(a_j)]$).



return the median of values in blue cells

2-wise hash function family

- 2-wise hash function definition:

- A family of hash function $H = \{h : U \rightarrow \{0, 1, \dots, w-1\}\}$ is **2-wise independent** if for any 2 distinct keys $x_i \neq x_j \in U$ and 2 hash values (not necessary distinct) $y_i, y_j \in \{0, 1, \dots, w-1\}$, we have

$$\Pr_h[h(x_i) = y_i \wedge h(x_j) = y_j] = 1/w^2$$

- On our CountSketch:

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

$$\Pr_h[h(a_i) = h(a_j)] = 1/w^2 * w = 1/w$$

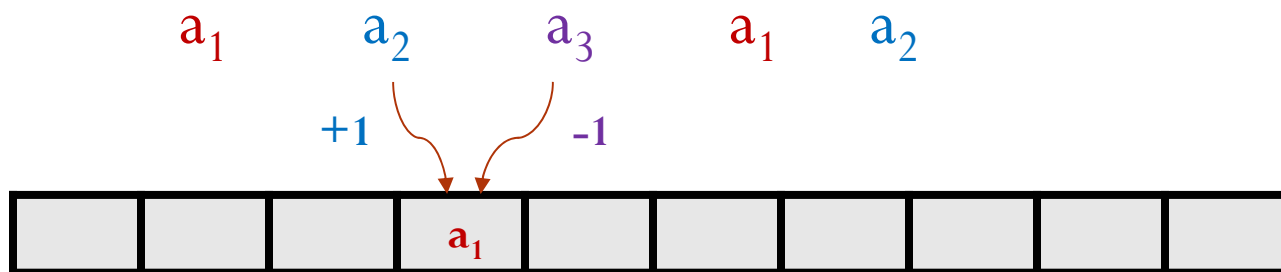
Analysis on 1 array of counters

- Notation:

- Stream of m items $\{a_1, \dots, a_m\}$ from the universe \mathbf{U} of size n .
- Frequency vector $\mathbf{f} = \{f_1, \dots, f_n\}$, $\|\mathbf{f}\|_2^2 = \sum_{i=1}^n f_i^2$

- Question:

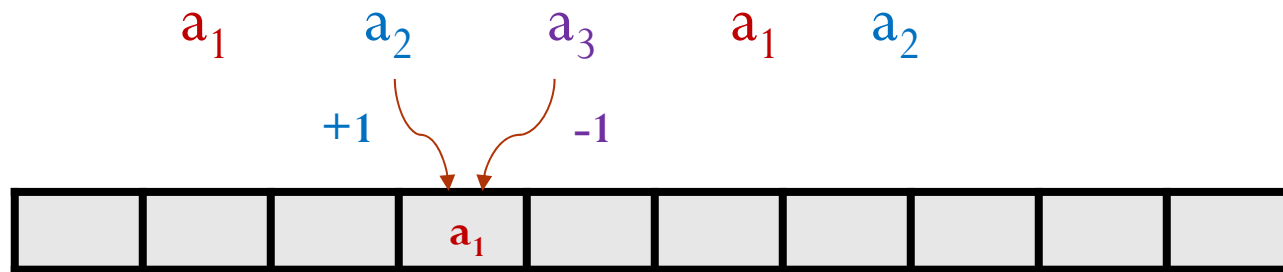
- Given a particular item a_1 , how many times $a_i \neq a_1$ collide by \mathbf{h} ?



Analysis on 1 array of counters

- Question:

- Given a particular item \mathbf{a}_1 , how many times $\mathbf{a}_i \neq \mathbf{a}_1$ collide by \mathbf{h} .



- Answer:

- Let \mathbf{X}_i be contribution of \mathbf{a}_i in the bucket $\mathbf{h}(\mathbf{a}_1)$.
- $$X_i = \begin{cases} f_i, & \text{which happens with prob. } 1/2w \\ 0, & \text{which happens with prob. } 1 - 1/w \\ -f_i, & \text{which happens with prob. } 1/2w \end{cases}$$
- Error source caused by the item \mathbf{a}_i : $\mathbf{E}[\mathbf{X}_i] = 0$

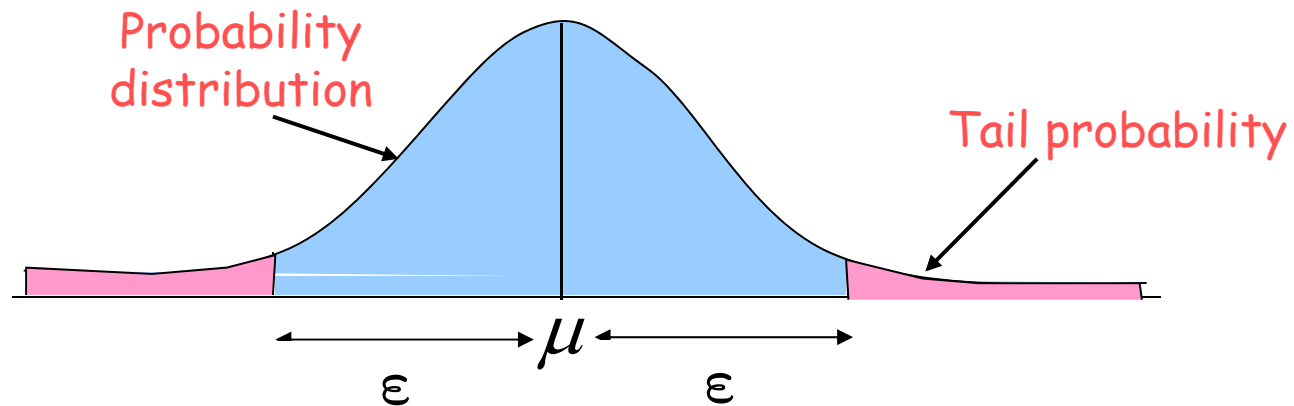
Analysis on 1 array of counters

- **Observation:** Our estimate is unbiased.
- **Question:** How large error do we have?
- **Analysis of variance of contribution of a particular item \mathbf{a}_i :**
 - Let \mathbf{X}_i be contribution of \mathbf{a}_i in the bucket $\mathbf{h}(\mathbf{a}_i)$.
 - $$X_i = \begin{cases} f_i, & \text{which happens with prob. } 1/2w \\ 0, & \text{which happens with prob. } 1 - 1/w \\ -f_i, & \text{which happens with prob. } 1/2w \end{cases}$$
 - Expectation: $\mathbf{E}[\mathbf{X}_i] = 0$.
 - Variance: $\mathbf{Var} [\mathbf{X}_i] = \mathbf{E}[\mathbf{X}_i^2] - (\mathbf{E}[\mathbf{X}_i])^2 = f_i^2/w$.

Analysis on 1 array of counters

- **Observation:** Our estimate is unbiased.
- **Question:** How large error do we have?
- **Analysis of variance of error:**
 - Let X_2, \dots, X_n be contributions of a_2, \dots, a_n in the bucket $h(a_1)$.
 - Let $Y = X_2 + \dots + X_n$ be the total contributions by other item $a_i \neq a_1$.
 - **Variance of error:**
$$\text{Var}[Y] = \text{Var}[X_2] + \dots + \text{Var}[X_n] = (f_2^2 + \dots + f_n^2)/w \leq \|f\|_2^2/w.$$

Basic tools: Tail inequality



- Chebyshev's inequality for $\mathbf{E}[Y] = \mu$ and $\mathbf{Var}[Y] = \sigma^2$:

$$\Pr [|Y - \mu| \geq \varepsilon] \leq \frac{\sigma^2}{\varepsilon^2} \text{ for any } \varepsilon > 0.$$

Analysis on 1 array of counters

- **Question:** How large error do we have?
- **Analysis of variance of error:**
 - Let X_2, \dots, X_n be contributions of a_2, \dots, a_n in the bucket $h(a_1)$.
 - Let $Y = X_2 + \dots + X_n$ be the total contributions by other item $a_i \neq a_1$.
 - Expectation: $E[Y] = 0$.
 - Variance: $\text{Var}[Y] = \text{Var}[X_2] + \dots + \text{Var}[X_n] = (f_2^2 + \dots + f_n^2)/w \leq \frac{\|f\|_2^2}{w}$
- **Chebyshev's inequality:**
 - $\Pr [|Y| \geq \epsilon \|f\|_2] \leq \frac{1}{w\epsilon^2} = 1/2$ if we choose $w = 2/\epsilon^2$.
 - The error is at most $\epsilon \|f\|_2$ with the probability $1/2$.

Analysis on d arrays of counters

- Boosting the accuracy:

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

- Analysis:

- $E[f'_1] = E[f'_2] = \dots = E[f'_d] = f_1$.
- Choose $d = \log(1/\delta)$ and apply Chernoff bound, we have

$$\Pr[|F_1 - f_1| \leq \epsilon \|f\|_2] \geq 1 - \delta$$

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

$$f_1 - \epsilon \|f\|_2 \leq F_1 \leq f_1 + \epsilon \|f\|_2$$

Homework

- Implement the CountSketch 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?
- What are the average errors of CountMin Sketch and CountSketch?