

## STREAMING ALGORITHMS

- Streaming Model of Computation
- Streaming Algorithms and DFA
- Stream: Motivation and Applications
- Synopsis: Sliding Window, Histogram, Wavelets
- Sampling from Stream: Reservoir Sampling
- Linear Sketch
- Count-Min Sketch
- AMS Sketch

IMDAD ULLAH KHAN

# Count-Min Sketch

## Count-min sketch

---

Count-min sketch stores frequencies of random groups of elements

▷ Cormode & Muthukrishnan (2004)

$$\mathcal{S} = \langle a_1, a_2, a_3, \dots, a_m \rangle \quad a_i \in [n]$$

$f_j$  : frequency of  $j$  in  $\mathcal{S}$        $\mathbf{F} = (f_1, f_2, \dots, f_n)$

- Cannot store frequency of every element  $j \in [n]$
- Store total frequency of random groups in  $[n]$  (elements in hash buckets)

---

**Algorithm** : Count-Min Sketch ( $k, \epsilon, \delta$ )

---

COUNT  $\leftarrow$  ZEROS( $k$ )

▷ sketch consists of  $k$  integers

Pick a random  $h : [n] \mapsto [k]$  from a 2-universal family  $\mathcal{H}$

On input  $a_i$

COUNT[ $h(a_i)$ ]  $\leftarrow$  COUNT[ $h(a_i)$ ] + 1

▷ increment count at index  $h(a_i)$

On query  $j$

▷ query:  $\mathbf{F}[j] = ?$

**return** COUNT[ $h(j)$ ]

---

# Count-min sketch

## Count-min sketch stores frequencies of random groups of elements

**Algorithm** : Count-Min Sketch ( $k, \epsilon, \delta$ )

COUNT  $\leftarrow$  ZEROS( $k$ )

▷ sketch consists of  $k$  integers

Pick a random  $h : [n] \mapsto [k]$  from a 2-universal family  $\mathcal{H}$

On input  $a_i$

COUNT[ $h(a_i)$ ]  $\leftarrow$  COUNT[ $h(a_i)$ ] + 1

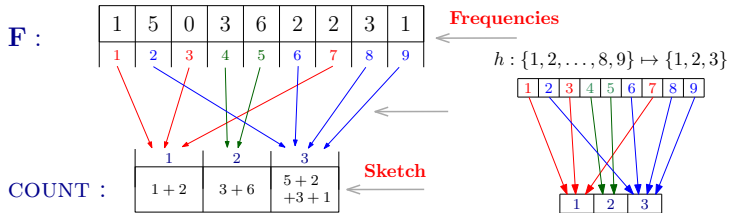
▷ increment count at index  $h(a_i)$

On query  $j$

return COUNT[ $h(j)$ ]

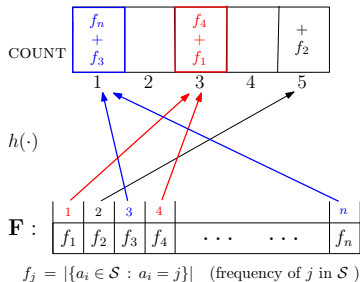
▷ query:  $\mathbf{F}[j] = ?$

$S$  : 2, 5, 6, 7, 8, 2, 1, 2, 7, 5, 5, 4, 2, 8, 8, 9, 5, 6, 4, 4, 2, 5, 5



## Count-min sketch

Count-min sketch stores frequencies of random groups of elements



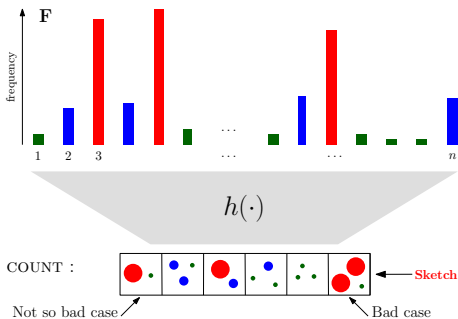
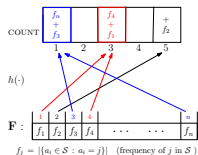
- $k = 2/\epsilon$
- Large  $k$  means better estimate (many smaller groups) but more space
- $\tilde{f}_j$ : estimate for  $f_j$  – output of algorithm

# Count-min sketch

Count-min sketch stores frequencies of random groups of elements

- $k = 2/\epsilon$
- Large  $k$  means better estimate but more space
- $\tilde{f}_j$ : estimate for  $f_j$  – output of algorithm

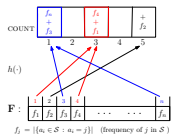
Quality on  $\tilde{f}_j$  :



# Count-min sketch

## Count-min sketch stores frequencies of random groups of elements

- $k = 2/\epsilon$
- Large  $k$  means better estimate but more space
- $\tilde{f}_j$ : estimate for  $f_j$  – output of algorithm



Quality on  $\tilde{f}_j$  :

- 1  $\tilde{f}_j \geq f_j$ 
  - Other elements that hash to  $h(j)$  contribute to  $\tilde{f}_j$
- 2  $Pr[\tilde{f}_j \leq f_j + \epsilon \|F\|_1] \geq 1/2$ 
  - $\tilde{f}_j = f_j + X_j$  ▷  $X_j$ : Excess in  $\tilde{f}_j$  (error)
  - $X_j = \sum_{i \in [n] \setminus j} f_i \cdot \mathbf{1}_{h(i)=h(j)}$  ▷  $\mathbf{1}_{condition}$  is indicator of condition

$$\mathbb{E}(X_j) = \mathbb{E}\left(\sum_{i \in [n] \setminus j} f_i \cdot \mathbf{1}_{h(i)=h(j)}\right) = \sum_{i \in [n] \setminus j} f_i \cdot \frac{1}{k} \leq \|F\|_1 \cdot \frac{\epsilon}{2}$$

- By Markov inequality we get the bound

## Count-min sketch

---

### Amplifying the probability of basic Count-Min Sketch

Keep  $t$  over-estimates,  $t = \log(1/\delta)$ ,  $k = 2/\epsilon$  and return their minimum

Unlikely that all  $t$  functions hash  $j$  with very frequent elements

---

**Algorithm** : Count-Min Sketch ( $k, \epsilon, \delta$ )

---

COUNT  $\leftarrow$  ZEROS( $t \times k$ ) ▷ sketch consists of  $t$  rows of  $k$  integers

Pick  $t$  random functions  $h_1, \dots, h_t : [n] \mapsto [k]$  from a 2-universal family

On input  $a_i$

**for**  $r = 1$  to  $t$  **do**

COUNT[ $r$ ][ $h_r(a_i)$ ]  $\leftarrow$  COUNT[ $r$ ][ $h_r(a_i)$ ] + 1  
▷ increment COUNT[ $r$ ] at index  $h_r(a_i)$

On query  $j$  ▷ query:  $\mathbf{F}[j] = ?$

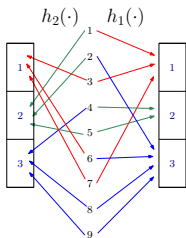
**return**  $\min_{1 \leq r \leq t}$  COUNT[ $r$ ][ $h_r(j)$ ]

---



# Count-min sketch

## Amplifying the probability of basic Count-Min Sketch



$\mathcal{S} : 2, 5, 6, 7, 8, 2, 1, 2, 7, 5, 5, 4, 2, 8, 8, 9, 5, 6, 4, 4, 2, 5, 5$

$h_1(\cdot)$

	1	2	3
COUNT :	0+1+2	3+6	5+2 +3+1
$h_2(\cdot)$	0+2+2	1+5+6	3+3+1

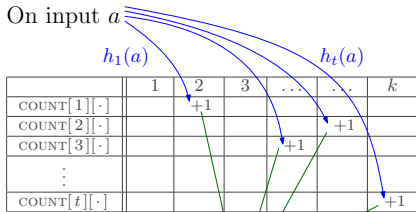
← Sketch

$F :$

1	2	3	4	5	6	7	8	9
1	5	0	3	6	2	2	3	1

← True Frequencies

On input  $a$



On query  $a$

$\text{MIN}_i \text{ COUNT}[i][h_i(a)]$

### Amplifying the probability of basic Count-Min Sketch

1  $\tilde{f}_j \geq f_j$

- For every  $r$ , other elements that hash to  $h_r(j)$  contribute to  $\tilde{f}_j$

2  $\Pr[\tilde{f}_j \leq f_j + \epsilon \|F\|_1] \geq 1 - \delta$

- $X_{jr}$  : contribution of other elements to  $\text{Count}[r][h_r(j)]$
- $\Pr[X_{jr} \geq \epsilon \|F\|_1] \leq 1/2$  for  $k = 2/\epsilon$
- The event  $\tilde{f}_j \geq f_j + \epsilon \|F\|_1$  is  $\forall 1 \leq r \leq t \quad X_{jr} \geq \epsilon \|F\|_1$
- $\Pr[\forall r \quad X_{jr} \geq \epsilon \|F\|_1] \leq (1/2)^t$
- $t = \log(1/\delta) \implies \Pr[\forall r \quad X_{jr} \geq \epsilon \|F\|_1] \leq (1/2)^{\log(1/\delta)} = \delta$

- Count-Min sketch is an  $(\epsilon \|F\|_1, \delta)$ -additive approximation algorithm
- Space required is  $k \cdot t$  integers =  $O(1/\epsilon \log(1/\delta) \log n)$  (plus constant)