

# Streaming: Sketching

---

Inge Li Gørtz

## Today

---

- Sketching
- CountMin sketch

# Sketching

## Sketching

---

- **Sketching.** create compact sketch/summary of data.
- **Example.** Durand and Flajolet 2003.
  - Condensed the whole Shakespeares' work
- Estimated number of distinct words: 30897 (correct answer is 28239, ie. relative error of 9.4%).
- **Composable.**
  - Data streams  $S_1$  and  $S_2$  with sketches  $sk(S_1)$  and  $sk(S_2)$
  - There exists an efficiently computable function  $f$  such that

```
ghfffghfghggggggghheehfhfhggghghhfgfffhhiigfhhffgfi ihfhhh  
igigighfgihfffghigihghigfhgeeghggghhghhfhidiigihighi ehhhfgg  
hfgighigffghdieghhggghhghhfiieffghghihifggffihgihfggighiif  
fjgfgjhjiifhjgehgghhfhjhiggghghihggghihighgiighghlghjgjjmfl
```

$$sk(S_1 \cup S_2) = f(sk(S_1), sk(S_2))$$

## Hashing

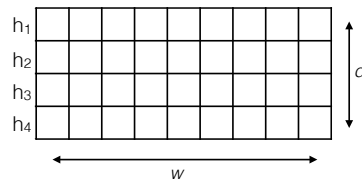
- **Pairwise independent hash function.** Let  $h : [n] \rightarrow [m]$ . For any  $x_1, x_2 \in [n]$  and  $y_1, y_2 \in [m]$  we have

$$\Pr[h(x_1) = y_1, h(x_2) = y_2] = \Pr[h(x_1) = y_1] \cdot \Pr[h(x_2) = y_2]$$

## CountMin Sketch

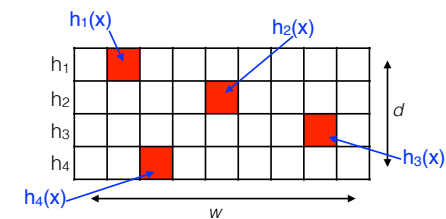
## CountMin Sketch

- Fixed array of counters of width  $w$  and depth  $d$ . Counters all initialized to be zero.
- Pairwise independent hash function for each row  $h_i : [n] \rightarrow [w]$ .
- When item  $x$  arrives increment counter  $h_i(x)$  of in all rows.



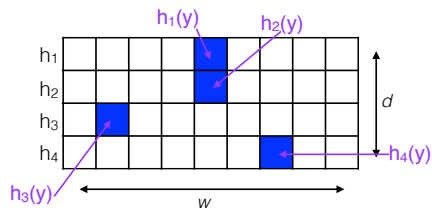
## CountMin Sketch

- Fixed array of counters of width  $w$  and depth  $d$ . Counters all initialized to be zero.
- Pairwise independent hash function for each row  $h_i : [n] \rightarrow [w]$ .
- When item  $x$  arrives increment counter  $h_i(x)$  of in all rows.



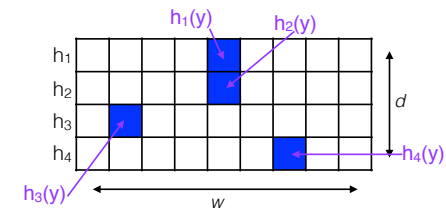
## CountMin Sketch

- Fixed array of counters of width  $w$  and depth  $d$ . Counters all initialized to be zero.
- Pairwise independent hash function for each row  $h_i : [n] \rightarrow [w]$ .
- When item  $x$  arrives increment counter  $h_i(x)$  of in all rows.
- Estimate frequency of  $y$ : return minimum of all entries  $y$  hash to.



## CountMin Sketch

- Fixed array of counters of width  $w$  and depth  $d$ . Counters all initialized to be zero.
- Pairwise independent hash function for each row  $h_i : [n] \rightarrow [w]$ .
- When item  $x$  arrives increment counter  $h_i(x)$  of in all rows.
- Estimate frequency of  $y$ : return minimum of all entries  $y$  hash to.



## CountMin Sketch

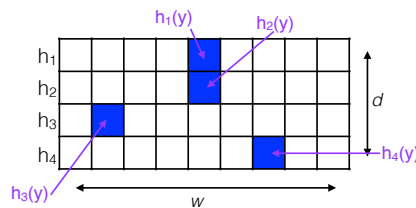
### Algorithm 1: CountMin

Initialize  $d$  independent hash functions  $h_j : [n] \rightarrow [w]$ .  
Set counter  $C[j, b] = 0$  for all  $j \in [n]$  and  $b \in [w]$ .

```

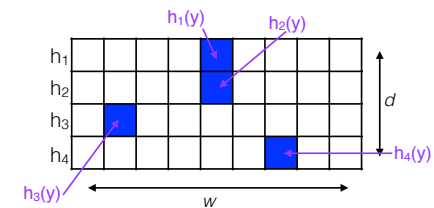
if  $Insert(x)$  then
  while  $Stream\ S\ not\ empty$  do
    for  $j = 1 \dots d$  do
       $C[j, h_j(x)] = +1$ 
    end
  end
else if  $Frequency(i)$  then
  return  $\hat{f}_i = \min_{j \in [d]} C(h_j(i))$ .
end
  
```

- The estimator  $\hat{f}_i$  has the following property:
  - $\hat{f}_i \geq f_i$
  - $\hat{f}_i \leq f_i + 2m/w$  with probability at least  $1 - (1/2)^d$



## CountMin Sketch: Analysis

- Use  $w = 2/\epsilon$  and  $d = \lg(1/\delta)$ .
- The estimator  $\hat{f}_i$  has the following property:
  - $\hat{f}_i \geq f_i$
  - $\hat{f}_i \leq f_i + \epsilon m$  with probability at least  $1 - \delta$
- **Space.**  $O(dw) = O(2 \lg(1/\delta)/\epsilon) = O(\lg(1/\delta)/\epsilon)$
- **Query and processing time.**  $O(d) = O(\lg(1/\delta))$



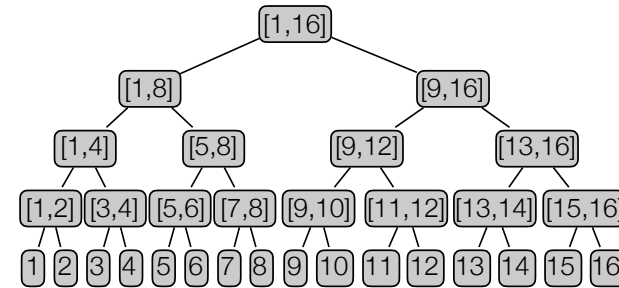
## Applications of CountMin Sketch

- We can use the CountMin Sketch to solve e.g.:
  - Heavy hitters:** List all heavy hitters (elements with frequency at least  $m/k$ ).
  - Range(a,b):** Return (an estimate of) the number of elements in the stream with value between a and b.

## Dyadic Intervals

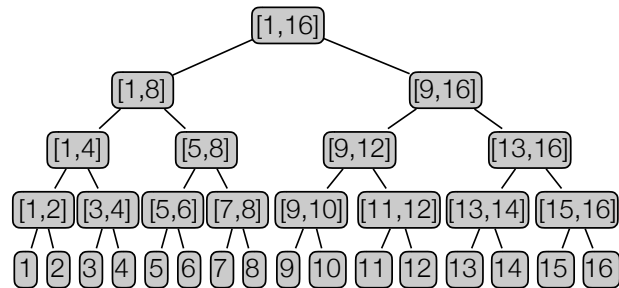
- Dyadic intervals.** Set of intervals:

$$\{[j\frac{n}{2^i} + 1, \dots, (j+1)\frac{n}{2^i}] \mid 0 \leq i \leq \lg n, 0 \leq j \leq 2^{i-1}\}$$



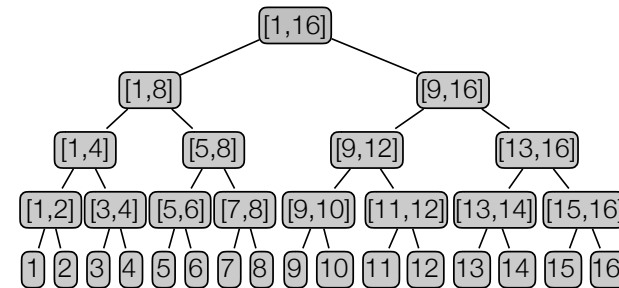
## Heavy Hitters

- Heavy Hitters.** Store a CountMin Sketch for each level in the tree of dyadic intervals
  - On a level: Treat all elements in same interval as the same element.



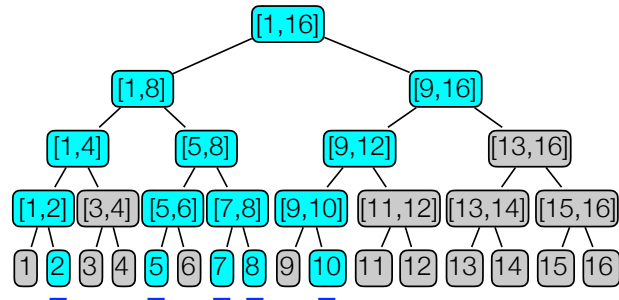
## Heavy Hitters

- Heavy Hitters.** Store a CountMin Sketch for each level in the tree of dyadic intervals
  - On a level: Treat all elements in same interval as the same element.



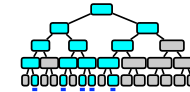
## Heavy Hitters

- **Heavy Hitters.** Store a CountMin Sketch for each level in the tree of dyadic intervals
  - On a level: Treat all elements in same interval as the same element.
- To find heavy hitters:
  - traverse tree from root.
  - only visit children with frequency  $\geq m/k$ .



## Heavy Hitters

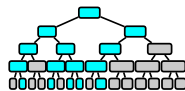
- **Heavy Hitters.** Store a CountMin Sketch for each level in the tree of dyadic intervals
  - On a level: Treat all elements in same interval as the same element.
- To find heavy hitters:
  - traverse tree from root.
  - only visit children with frequency  $\geq m/k$ .



## Heavy Hitters

- **Heavy Hitters.** Store a CountMin Sketch for each level in the tree of dyadic intervals
  - On a level: Treat all elements in same interval as the same element.
- To find heavy hitters:
  - traverse tree from root.
  - only visit children with frequency  $\geq m/k$ .
- **Analysis.**
  - **Time.**
    - Number of intervals queried:  $O(k \lg n)$ .
    - Query time:  $O(k \lg n \cdot \lg(1/\delta))$
  - **Space.**

$$O\left(\lg n \cdot \frac{1}{\epsilon} \lg\left(\frac{1}{\delta}\right)\right)$$



## Count Sketch

### Algorithm 2: CountSketch

```

Initialize  $d$  independent hash functions  $h_j : [n] \rightarrow [w]$ .
Initialize  $d$  independent hash functions  $s_j : [n] \rightarrow \{\pm 1\}$ .
Set counter  $C[j, b] = 0$  for all  $j \in [d]$  and  $b \in [w]$ .
if  $Insert(x)$  then
    while  $Stream S$  not empty do
        for  $j = 1 \dots d$  do
             $C[j, h_j(x)] =+ s_j(i)$ 
        end
    end
else if  $Frequency(i)$  then
     $\hat{f}_{ij} = C(h_j(i)) \cdot s_j(i)$ 
    return  $\hat{f}_{ij} = \text{median}_{j \in [d]} \hat{f}_{ij}$ 
end
    
```

	Space	Error
Count-Min	$O\left(\frac{1}{\epsilon} \log n\right)$	$\epsilon F_1$ (one-sided)
Count-Sketch	$O\left(\frac{1}{\epsilon^2} \log n\right)$	$\epsilon \sqrt{F_2}$ (two-sided)