

# Topics in TCS

---

**Frequency estimation via sketching**

---

Raphaël Clifford

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

The MISRA-GRIES algorithm is one-pass and runs in  $O(k(\log m + \log n))$  bits of space and  $O(m \log n)$  time.

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

The MISRA-GRIES algorithm is one-pass and runs in  $O(k(\log m + \log n))$  bits of space and  $O(m \log n)$  time.

It is deterministic (good✓) but only works in the cash register model.

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

The MISRA-GRIES algorithm is one-pass and runs in  $O(k(\log m + \log n))$  bits of space and  $O(m \log n)$  time.

It is deterministic (good✓) but only works in the cash register model.

We will change the definition to ask for an estimate of the frequency of occurrence for any token queried.

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

The MISRA-GRIES algorithm is one-pass and runs in  $O(k(\log m + \log n))$  bits of space and  $O(m \log n)$  time.

It is deterministic (good✓) but only works in the cash register model.

We will change the definition to ask for an estimate of the frequency of occurrence for any token queried.

We will introduce our first randomised *sketching* algorithms.

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

The MISRA-GRIES algorithm is one-pass and runs in  $O(k(\log m + \log n))$  bits of space and  $O(m \log n)$  time.

It is deterministic (good✓) but only works in the cash register model.

We will change the definition to ask for an estimate of the frequency of occurrence for any token queried.

We will introduce our first randomised *sketching* algorithms.

Our sketches will be *linear* which will mean we can extend them to the *turnstile model*. We can also combine them easily by adding.

## Frequent items via sketching

We return to the problem of finding frequent items. Our previous definition was: Given a parameter  $k$ , find the set of symbols with frequency greater than  $m/k$ .

The MISRA-GRIES algorithm is one-pass and runs in  $O(k(\log m + \log n))$  bits of space and  $O(m \log n)$  time.

It is deterministic (good✓) but only works in the cash register model.

We will change the definition to ask for an estimate of the frequency of occurrence for any token queried.

We will introduce our first randomised *sketching* algorithms.

Our sketches will be *linear* which will mean we can extend them to the *turnstile model*. We can also combine them easily by adding.

They will give us an estimate of the frequency for *every* token.



## COUNTSKETCH

The sketch is a 2D-array  $C$  with  $t$  rows and  $k$  columns. All hash functions are chosen from a pairwise independent family.

# COUNTSKETCH

The sketch is a 2D-array  $C$  with  $t$  rows and  $k$  columns. All hash functions are chosen from a pairwise independent family.

```
stream  $\langle a_1, \dots, a_m \rangle, a_i \in [n]$   
initialise  $C[1 \dots t][1 \dots k] = 0$   
choose hash functions  $h_1, \dots, h_t : [n] \rightarrow [k]$   
choose hash function  $g_1, \dots, g_t : [n] \rightarrow \{-1, 1\}$ 
```





```
COUNTSKETCH( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_j g_j(a_i)$   
  
return  $\hat{f}_{a_i} = \text{median}\{g_j(a_i) C[j, h_j(a_i)]\}$ 
```

$c_i$  is the number of instances of  $a_i$ . In the turnstile model this can be either positive or negative.

# COUNTSKETCH - worked example



	1	2	3
$h_1$			
$h_2$			




	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```
COUNTSKETCH( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_j g_j(a_i)$ 
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$		+	
$h_2$	+		



	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```
COUNTSKETCH( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$		+	-
$h_2$	+	+	

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```
COUNTSKETCH( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_j g_j(a_i)$ 
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$		++	-
$h_2$	++	+	

	$h_1, g_1$	$h_2, g_2$
●	2, +	1, +
●	3, -	2, +
●	1, +	3, -
●	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$		++-	-
$h_2$	++	+	+

	$h_1, g_1$	$h_2, g_2$
●	2, +	1, +
●	3, -	2, +
●	1, +	3, -
●	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$		++-+	-
$h_2$	+++	+	+

	$h_1, g_1$	$h_2, g_2$
●	2, +	1, +
●	3, -	2, +
●	1, +	3, -
●	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```



# COUNTSKETCH - worked example



	1	2	3
$h_1$		++-+	--
$h_2$	+++	++	+

	$h_1, g_1$	$h_2, g_2$
●	2, +	1, +
●	3, -	2, +
●	1, +	3, -
●	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	+	++-+	--
$h_2$	++++	++	+-

	$h_1, g_1$	$h_2, g_2$
●	2, +	1, +
●	3, -	2, +
●	1, +	3, -
●	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	+	++-++	--
$h_2$	+++++	++	+-

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	+	++-++ +	--
$h_2$	+++++	++	+-

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	+	++-++ +	---
$h_2$	+++++	+++	+-

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +	---
$h_2$	+++++	+++	+--

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_j g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +-	---
$h_2$	+++++	+++	+---+

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

```

COUNTSKETCH( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i g_j(a_i)$ 
    
```

# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +-	---
$h_2$	+++++	+++	+---+

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

return  $\hat{f}_{a_i} = \text{median}\{g_j(a_i)C[j, h_j(a_i)]\}$



# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +-	---
$h_2$	+++++	+++	+---+

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

return  $\hat{f}_{a_i} = \text{median}\{g_j(a_i)C[j, h_j(a_i)]\}$

$$\hat{f}_{\text{green}} = \text{median}(g_1(\text{green})C[1, h_1(\text{green})], g_2(\text{green})C[2, h_2(\text{green})]) = \text{median}(1 \cdot 3, 1 \cdot 5)$$

# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +-	---
$h_2$	+++++	+++	+---+

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

return  $\hat{f}_{a_i} = \text{median}\{g_j(a_i)C[j, h_j(a_i)]\}$

$$\hat{f}_{\text{green}} = \text{median}(g_1(\text{green})C[1, h_1(\text{green})], g_2(\text{green})C[2, h_2(\text{green})]) = \text{median}(1 \cdot 3, 1 \cdot 5)$$

$$\hat{f}_{\text{cyan}} = \text{median}(g_1(\text{cyan})C[1, h_1(\text{cyan})], g_2(\text{cyan})C[2, h_2(\text{cyan})]) = \text{median}(-1 \cdot -3, 1 \cdot 3)$$

# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +-	---
$h_2$	+++++	+++	+---+

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

return  $\hat{f}_{a_i} = \text{median}\{g_j(a_i)C[j, h_j(a_i)]\}$

$$\hat{f}_{\text{green}} = \text{median}(g_1(\text{green})C[1, h_1(\text{green})], g_2(\text{green})C[2, h_2(\text{green})]) = \text{median}(1 \cdot 3, 1 \cdot 5)$$

$$\hat{f}_{\text{cyan}} = \text{median}(g_1(\text{cyan})C[1, h_1(\text{cyan})], g_2(\text{cyan})C[2, h_2(\text{cyan})]) = \text{median}(-1 \cdot -3, 1 \cdot 3)$$

$$\hat{f}_{\text{purple}} = \text{median}(g_1(\text{purple})C[1, h_1(\text{purple})], g_2(\text{purple})C[2, h_2(\text{purple})]) = \text{median}(1 \cdot 2, -1 \cdot 0)$$

# COUNTSKETCH - worked example



	1	2	3
$h_1$	++	++-++ +-	---
$h_2$	+++++	+++	+---+

	$h_1, g_1$	$h_2, g_2$
	2, +	1, +
	3, -	2, +
	1, +	3, -
	2, -	3, +

return  $\hat{f}_{a_i} = \text{median}\{g_j(a_i)C[j, h_j(a_i)]\}$

$$\hat{f}_{\text{green}} = \text{median}(g_1(\text{green})C[1, h_1(\text{green})], g_2(\text{green})C[2, h_2(\text{green})]) = \text{median}(1 \cdot 3, 1 \cdot 5)$$

$$\hat{f}_{\text{cyan}} = \text{median}(g_1(\text{cyan})C[1, h_1(\text{cyan})], g_2(\text{cyan})C[2, h_2(\text{cyan})]) = \text{median}(-1 \cdot -3, 1 \cdot 3)$$

$$\hat{f}_{\text{purple}} = \text{median}(g_1(\text{purple})C[1, h_1(\text{purple})], g_2(\text{purple})C[2, h_2(\text{purple})]) = \text{median}(1 \cdot 2, -1 \cdot 0)$$

$$\hat{f}_{\text{red}} = \text{median}(g_1(\text{red})C[1, h_1(\text{red})], g_2(\text{red})C[2, h_2(\text{red})]) = \text{median}(-1 \cdot 3, 1 \cdot 0)$$

## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

Let  $X = \hat{f}_a$  be the output for query  $a$ .

## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

Let  $X = \hat{f}_a$  be the output for query  $a$ .

For each token  $j$ , define indicator r.v.  $Y_j = 1$  iff  $h(j) = h(a)$ .

## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

Let  $X = \hat{f}_a$  be the output for query  $a$ .

For each token  $j$ , define indicator r.v.  $Y_j = 1$  iff  $h(j) = h(a)$ .

Token  $j$  contributes  $f_j \cdot g(j)$  to  $C[h(a)]$  iff  $h(j) = h(a)$ .



## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

Let  $X = \hat{f}_a$  be the output for query  $a$ .

For each token  $j$ , define indicator r.v.  $Y_j = 1$  iff  $h(j) = h(a)$ .

Token  $j$  contributes  $f_j \cdot g(j)$  to  $C[h(a)]$  iff  $h(j) = h(a)$ .

Therefore

$$X = g(a) \sum_{j=1}^n f_j g(j) Y_j = f_a + \sum_{j \in [n] \setminus \{a\}} f_j g(a) g(j) Y_j$$

## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

Let  $X = \hat{f}_a$  be the output for query  $a$ .

For each token  $j$ , define indicator r.v.  $Y_j = 1$  iff  $h(j) = h(a)$ .

Token  $j$  contributes  $f_j \cdot g(j)$  to  $C[h(a)]$  iff  $h(j) = h(a)$ .

Therefore

$$X = g(a) \sum_{j=1}^n f_j g(j) Y_j = f_a + \sum_{j \in [n] \setminus \{a\}} f_j g(a) g(j) Y_j$$

As  $g$  and  $h$  are independent and  $g$  is from a pairwise independent family,

$$\mathbb{E}[g(a)g(j)Y_j] = \mathbb{E}(g(a)) \cdot \mathbb{E}(g(j)) \cdot \mathbb{E}(Y_j) = 0 \cdot 0 \cdot \mathbb{E}(Y_j) = 0$$

## COUNTSKETCH - Analysis I

To start, let us look just at an arbitrary row of  $C$ . We will show that for each row COUNTSKETCH gives an unbiased estimate. Define  $C[x] = C[1, x]$ .

Let  $X = \hat{f}_a$  be the output for query  $a$ .

For each token  $j$ , define indicator r.v.  $Y_j = 1$  iff  $h(j) = h(a)$ .

Token  $j$  contributes  $f_j \cdot g(j)$  to  $C[h(a)]$  iff  $h(j) = h(a)$ .

Therefore

$$X = g(a) \sum_{j=1}^n f_j g(j) Y_j = f_a + \sum_{j \in [n] \setminus \{a\}} f_j g(a) g(j) Y_j$$

As  $g$  and  $h$  are independent and  $g$  is from a pairwise independent family,

$$\mathbb{E}[g(a)g(j)Y_j] = \mathbb{E}(g(a)) \cdot \mathbb{E}(g(j)) \cdot \mathbb{E}(Y_j) = 0 \cdot 0 \cdot \mathbb{E}(Y_j) = 0$$

By linearity of expectation

$$\mathbb{E}(X) = f_a + \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a)g(j)Y_j] = f_a$$

## COUNTSKETCH - Analysis IIa

We will now derive the variance of our estimator  $X = \hat{f}$ . Recall  $Y_j = 1$  iff  $h(j) = h(a)$ .

## COUNTSKETCH - Analysis IIa

We will now derive the variance of our estimator  $X = \hat{f}$ . Recall  $Y_j = 1$  iff  $h(j) = h(a)$ .

$$\text{var}(X) = 0 + \text{var} \left[ g(a) \sum_{j \in [n] \setminus \{a\}} f_j \cdot g(j) Y_j \right]$$

## COUNTSKETCH - Analysis IIa

We will now derive the variance of our estimator  $X = \hat{f}$ . Recall  $Y_j = 1$  iff  $h(j) = h(a)$ .

$$\begin{aligned}\text{var}(X) &= 0 + \text{var} \left[ g(a) \sum_{j \in [n] \setminus \{a\}} f_j \cdot g(j) Y_j \right] \\ &= \mathbb{E} \left[ g(a)^2 \sum_{j \in [n] \setminus \{a\}} f_j^2 Y_j^2 + \sum_{\substack{j \in [n] \setminus \{a\} \\ i \neq j}} f_i f_j g(i) g(j) Y_i Y_j \right] - \\ &\quad \left[ \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a) g(j) Y_j] \right]^2\end{aligned}$$

## COUNTSKETCH - Analysis IIa

We will now derive the variance of our estimator  $X = \hat{f}$ . Recall  $Y_j = 1$  iff  $h(j) = h(a)$ .

$$\begin{aligned}\text{var}(X) &= 0 + \text{var} \left[ g(a) \sum_{j \in [n] \setminus \{a\}} f_j \cdot g(j) Y_j \right] \\ &= \mathbb{E} \left[ g(a)^2 \sum_{j \in [n] \setminus \{a\}} f_j^2 Y_j^2 + \sum_{\substack{j \in [n] \setminus \{a\} \\ i \neq j}} f_i f_j g(i) g(j) Y_i Y_j \right] - \\ &\quad \left[ \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a) g(j) Y_j] \right]^2\end{aligned}$$

We will need two facts to simplify these terms.

## COUNTSKETCH - Analysis IIb

$$\text{var}(X) = \mathbb{E} \left[ g(a)^2 \sum_{j \in [n] \setminus \{a\}} f_j^2 Y_j^2 + \sum_{\substack{j \in [n] \setminus \{a\} \\ i \neq j}} f_i f_j g(i) g(j) Y_i Y_j \right] - \left[ \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a) g(j) Y_j] \right]^2$$



## COUNTSKETCH - Analysis IIb

$$\text{var}(X) = \mathbb{E} \left[ g(a)^2 \sum_{j \in [n] \setminus \{a\}} f_j^2 Y_j^2 + \sum_{\substack{j \in [n] \setminus \{a\} \\ i \neq j}} f_i f_j g(i) g(j) Y_i Y_j \right] - \left[ \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a) g(j) Y_j] \right]^2$$

Now, the two facts:

1.  $\mathbb{E}(Y_j^2) = \mathbb{E}(Y_j) = \Pr(h(j) = h(a)) = \frac{1}{k}$ .
2.  $\mathbb{E}(g(i)g(j)Y_iY_j) = \mathbb{E}(g(i)) \cdot \mathbb{E}(g(j)) \cdot \mathbb{E}(Y_iY_j) = 0 \cdot 0 \cdot \mathbb{E}(Y_iY_j) = 0$

## COUNTSKETCH - Analysis IIb

$$\text{var}(X) = \mathbb{E} \left[ g(a)^2 \sum_{j \in [n] \setminus \{a\}} f_j^2 Y_j^2 + \sum_{\substack{j \in [n] \setminus \{a\} \\ i \neq j}} f_i f_j g(i) g(j) Y_i Y_j \right] - \left[ \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a) g(j) Y_j] \right]^2$$

Now, the two facts:

1.  $\mathbb{E}(Y_j^2) = \mathbb{E}(Y_j) = \Pr(h(j) = h(a)) = \frac{1}{k}$ .
2.  $\mathbb{E}(g(i)g(j)Y_iY_j) = \mathbb{E}(g(i)) \cdot \mathbb{E}(g(j)) \cdot \mathbb{E}(Y_iY_j) = 0 \cdot 0 \cdot \mathbb{E}(Y_iY_j) = 0$

Therefore,

$$\text{var}(X) = \sum_{j \in [n] \setminus \{a\}} \frac{f_j^2}{k} + 0 - 0$$

## COUNTSKETCH - Analysis IIb

$$\text{var}(X) = \mathbb{E} \left[ g(a)^2 \sum_{j \in [n] \setminus \{a\}} f_j^2 Y_j^2 + \sum_{\substack{j \in [n] \setminus \{a\} \\ i \neq j}} f_i f_j g(i) g(j) Y_i Y_j \right] - \left[ \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}[g(a) g(j) Y_j] \right]^2$$

Now, the two facts:

1.  $\mathbb{E}(Y_j^2) = \mathbb{E}(Y_j) = \Pr(h(j) = h(a)) = \frac{1}{k}$ .
2.  $\mathbb{E}(g(i)g(j)Y_iY_j) = \mathbb{E}(g(i)) \cdot \mathbb{E}(g(j)) \cdot \mathbb{E}(Y_iY_j) = 0 \cdot 0 \cdot \mathbb{E}(Y_iY_j) = 0$

Therefore,

$$\begin{aligned} \text{var}(X) &= \sum_{j \in [n] \setminus \{a\}} \frac{f_j^2}{k} + 0 - 0 \\ &= \frac{\|\mathbf{f}\|_2^2 - f_a^2}{k} \quad \text{where } \mathbf{f} \text{ is the array of frequencies} \end{aligned}$$

## COUNTSKETCH - Analysis III

Using the variance  $\text{var}(X) = \frac{\|f\|_2^2 - f_a^2}{k}$  we can apply Chebyshev.

## COUNTSKETCH - Analysis III

Using the variance  $\text{var}(X) = \frac{\|\mathbf{f}\|_2^2 - f_a^2}{k}$  we can apply Chebyshev.

$$\begin{aligned}\Pr(|\hat{f}_a - f_a| \geq \epsilon \sqrt{\|\mathbf{f}\|_2^2 - f_a^2}) &= \Pr(|X - \mathbb{E}(X)| \geq \epsilon \sqrt{\|\mathbf{f}\|_2^2 - f_a^2}) \\ &\leq \frac{\text{var}(X)}{\epsilon^2(\|\mathbf{f}\|_2^2 - f_a^2)} \\ &= \frac{1}{k\epsilon^2} \\ &= \frac{1}{3} \quad (\text{set } k = 3/\epsilon^2)\end{aligned}$$

## COUNTSKETCH - Analysis III

Using the variance  $\text{var}(X) = \frac{\|\mathbf{f}\|_2^2 - f_a^2}{k}$  we can apply Chebyshev.

$$\begin{aligned}\Pr(|\hat{f}_a - f_a| \geq \epsilon \sqrt{\|\mathbf{f}\|_2^2 - f_a^2}) &= \Pr(|X - \mathbb{E}(X)| \geq \epsilon \sqrt{\|\mathbf{f}\|_2^2 - f_a^2}) \\ &\leq \frac{\text{var}(X)}{\epsilon^2(\|\mathbf{f}\|_2^2 - f_a^2)} \\ &= \frac{1}{k\epsilon^2} \\ &= \frac{1}{3} \quad (\text{set } k = 3/\epsilon^2)\end{aligned}$$

Using the notation  $\mathbf{f}_{-j}$  for  $\mathbf{f}$  with the  $j$ th element dropped,  
 $\|\mathbf{f}_{-j}\|_2^2 = \|\mathbf{f}\|_2^2 - f_j^2$ .

## COUNTSKETCH - Analysis III

Using the variance  $\text{var}(X) = \frac{\|\mathbf{f}\|_2^2 - f_a^2}{k}$  we can apply Chebyshev.

$$\begin{aligned}\Pr(|\hat{f}_a - f_a| \geq \epsilon \sqrt{\|\mathbf{f}\|_2^2 - f_a^2}) &= \Pr(|X - \mathbb{E}(X)| \geq \epsilon \sqrt{\|\mathbf{f}\|_2^2 - f_a^2}) \\ &\leq \frac{\text{var}(X)}{\epsilon^2(\|\mathbf{f}\|_2^2 - f_a^2)} \\ &= \frac{1}{k\epsilon^2} \\ &= \frac{1}{3} \quad (\text{set } k = 3/\epsilon^2)\end{aligned}$$

Using the notation  $\mathbf{f}_{-j}$  for  $\mathbf{f}$  with the  $j$ th element dropped,  $\|\mathbf{f}_{-j}\|_2^2 = \|\mathbf{f}\|_2^2 - f_j^2$ . And so,

$$\Pr(|\hat{f}_a - f_a| \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \frac{1}{3}$$

## COUNTSKETCH - Analysis IV

So how good is our sketch that takes the median?



## COUNTSKETCH - Analysis IV

So how good is our sketch that takes the median?

We take the median of  $|\hat{f}_a - f_a|$  for  $t$  different independent runs. If this is at least  $\epsilon \|\mathbf{f}_{-a}\|_2$  then at least  $t/2$  iterations are that big.

## COUNTSKETCH - Analysis IV

So how good is our sketch that takes the median?

We take the median of  $|\hat{f}_a - f_a|$  for  $t$  different independent runs. If this is at least  $\epsilon \|\mathbf{f}_{-a}\|_2$  then at least  $t/2$  iterations are that big.

We show that this is exponentially unlikely to happen as a function of the number of iterations,  $t$ .

## COUNTSKETCH - Analysis IV

So how good is our sketch that takes the median?

We take the median of  $|\hat{f}_a - f_a|$  for  $t$  different independent runs. If this is at least  $\epsilon \|\mathbf{f}_{-a}\|_2$  then at least  $t/2$  iterations are that big.

We show that this is exponentially unlikely to happen as a function of the number of iterations,  $t$ .

For the  $i$ th iteration, let  $Z_i = 1$  if  $|\hat{f}_a - f_a| \geq \epsilon \|\mathbf{f}_{-a}\|_2$  and 0 otherwise.

## COUNTSKETCH - Analysis IV

So how good is our sketch that takes the median?

We take the median of  $|\hat{f}_a - f_a|$  for  $t$  different independent runs. If this is at least  $\epsilon \|\mathbf{f}_{-a}\|_2$  then at least  $t/2$  iterations are that big.

We show that this is exponentially unlikely to happen as a function of the number of iterations,  $t$ .

For the  $i$ th iteration, let  $Z_i = 1$  if  $|\hat{f}_a - f_a| \geq \epsilon \|\mathbf{f}_{-a}\|_2$  and 0 otherwise.

Using Chernoff's bound with  $\mu = t/3$

$$\Pr\left(\sum_{i=1}^t Z_i \geq (1 + \delta)\mu\right) \leq \exp(-\delta^2\mu/3) = \exp(-\delta^2 t/9)$$

$$\Pr\left(\sum_{i=1}^t Z_i \geq (1 + 1/2)\mu\right) \leq \exp(-(1/2)^2 t/9) = \exp(-t/36)$$

## COUNTSKETCH - Analysis IV

So how good is our sketch that takes the median?

We take the median of  $|\hat{f}_a - f_a|$  for  $t$  different independent runs. If this is at least  $\epsilon \|\mathbf{f}_{-a}\|_2$  then at least  $t/2$  iterations are that big.

We show that this is exponentially unlikely to happen as a function of the number of iterations,  $t$ .

For the  $i$ th iteration, let  $Z_i = 1$  if  $|\hat{f}_a - f_a| \geq \epsilon \|\mathbf{f}_{-a}\|_2$  and 0 otherwise.

Using Chernoff's bound with  $\mu = t/3$

$$\Pr\left(\sum_{i=1}^t Z_i \geq (1 + \delta)\mu\right) \leq \exp(-\delta^2\mu/3) = \exp(-\delta^2 t/9)$$

$$\Pr\left(\sum_{i=1}^t Z_i \geq (1 + 1/2)\mu\right) \leq \exp(-(1/2)^2 t/9) = \exp(-t/36)$$

For an arbitrary token  $a$ , the probability of being further than  $\epsilon \|\mathbf{f}_{-a}\|_2$  from the correct frequency is at most  $\exp(-t/36)$ .

## COUNTSKETCH - Space/Time

We need  $O(\log m)$  bits per counter in our sketch. There are  $tk$  counters.

## COUNTSKETCH - Space/Time

We need  $O(\log m)$  bits per counter in our sketch. There are  $tk$  counters.

We store  $t$  pairwise independent hash functions making  $O(t \log n)$  bits.

## COUNTSKETCH - Space/Time

We need  $O(\log m)$  bits per counter in our sketch. There are  $tk$  counters.

We store  $t$  pairwise independent hash functions making  $O(t \log n)$  bits.

Overall space is therefore  $O(t \log n + tk \log m)$  bits.



## COUNTSKETCH - Space/Time

We need  $O(\log m)$  bits per counter in our sketch. There are  $tk$  counters.

We store  $t$  pairwise independent hash functions making  $O(t \log n)$  bits.

Overall space is therefore  $O(t \log n + tk \log m)$  bits.

With  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \ln 1/\delta \rceil$ , this equals

$$O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \cdot (\log m + \log n)\right) \text{ bits}$$

## COUNTSKETCH - Space/Time

We need  $O(\log m)$  bits per counter in our sketch. There are  $tk$  counters.

We store  $t$  pairwise independent hash functions making  $O(t \log n)$  bits.

Overall space is therefore  $O(t \log n + tk \log m)$  bits.

With  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \ln 1/\delta \rceil$ , this equals

$$O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \cdot (\log m + \log n)\right) \text{ bits}$$

Running time: one-pass and  $O(t)$  time per token.

## COUNTSKETCH summary

COUNTSKETCH is a one-pass randomised algorithm to estimate the frequency of the tokens in a stream.

## COUNTSKETCH summary

COUNTSKETCH is a one-pass randomised algorithm to estimate the frequency of the tokens in a stream.

Once  $\epsilon$  and  $\delta$  are decided we can set  $t$  and  $k$  accordingly.

## COUNTSKETCH summary

COUNTSKETCH is a one-pass randomised algorithm to estimate the frequency of the tokens in a stream.

Once  $\epsilon$  and  $\delta$  are decided we can set  $t$  and  $k$  accordingly.

The running time is  $O(t)$  time per token.

## COUNTSKETCH summary

COUNTSKETCH is a one-pass randomised algorithm to estimate the frequency of the tokens in a stream.

Once  $\epsilon$  and  $\delta$  are decided we can set  $t$  and  $k$  accordingly.

The running time is  $O(t)$  time per token.

The space usage is  $O(t \log n + tk \log m)$  bits.

## COUNTSKETCH summary

COUNTSKETCH is a one-pass randomised algorithm to estimate the frequency of the tokens in a stream.

Once  $\epsilon$  and  $\delta$  are decided we can set  $t$  and  $k$  accordingly.

The running time is  $O(t)$  time per token.

The space usage is  $O(t \log n + tk \log m)$  bits.

Assuming we set  $k = 3/\epsilon^2$ , for an arbitrary token  $a$ , the probability that COUNTSKETCH's estimate is further than  $\epsilon \|\mathbf{f}_a\|_2$  from the correct frequency is at most  $\exp(-t/36)$ .

## COUNT-MIN sketch

The sketch is a 2D-array  $C$  with  $t$  rows and  $k$  columns. All hash functions are chosen from a pairwise independent family.



## COUNT-MIN sketch

The sketch is a 2D-array  $C$  with  $t$  rows and  $k$  columns. All hash functions are chosen from a pairwise independent family.

```
stream  $\langle a_1, \dots, a_m \rangle, a_i \in [n]$   
initialise  $C[1 \dots t][1 \dots k] = 0$   
choose hash functions  $h_1, \dots, h_t : [n] \rightarrow [k]$ 
```

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

```
return  $\hat{f}_a = \min_{1 \leq i \leq t} C[i, h_i(a)]$ 
```

$c_i$  is the number of instances of  $a_i$ . In the turnstile model this can be either positive or negative.

# COUNT-MIN - worked example






	1	2	3
$h_1$			
$h_2$			
$h_3$			

	$h_1$	$h_2$	$h_3$
●	1	2	3
●	2	1	1
●	1	1	1
●	3	3	2

COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$

## COUNT-MIN - worked example









	1	2	3
$h_1$			
$h_2$			
$h_3$			

	$h_1$	$h_2$	$h_3$
	1	2	3
	2	1	1
	1	1	1
	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

## COUNT-MIN - worked example

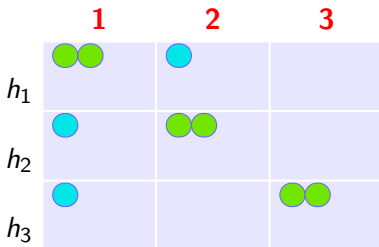


	1	2	3
$h_1$			
$h_2$			
$h_3$			

	$h_1$	$h_2$	$h_3$
	1	2	3
	2	1	1
	1	1	1
	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

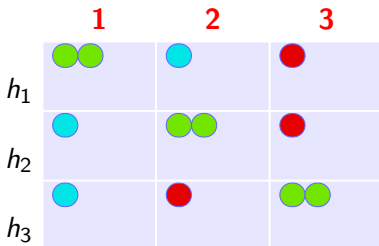
## COUNT-MIN - worked example



	$h_1$	$h_2$	$h_3$
Green	1	2	3
Cyan	2	1	1
Purple	1	1	1
Red	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

## COUNT-MIN - worked example


















	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

## COUNT-MIN - worked example



	1	2	3
$h_1$	  		
$h_2$		  	
$h_3$			  

	$h_1$	$h_2$	$h_3$
	1	2	3
	2	1	1
	1	1	1
	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

## COUNT-MIN - worked example



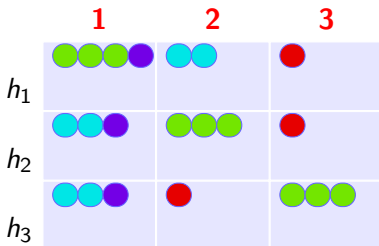
	1	2	3
$h_1$			
$h_2$			
$h_3$			

	$h_1$	$h_2$	$h_3$
	1	2	3
	2	1	1
	1	1	1
	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```



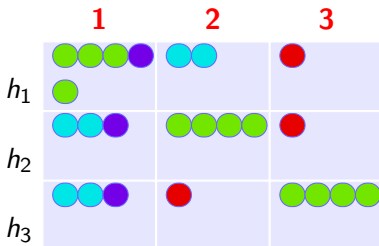
## COUNT-MIN - worked example



	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

```
COUNT-MIN( $a_i$ )  
for each  $j \in [t]$   
     $C[j, h_j(a_i)] += c_i$ 
```

# COUNT-MIN - worked example

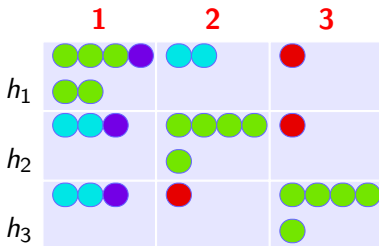


	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

```

COUNT-MIN( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i$ 
    
```

# COUNT-MIN - worked example

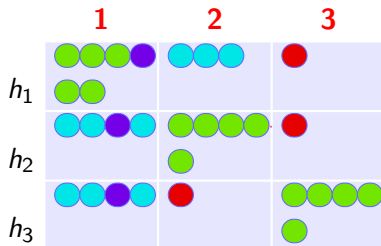


	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

```

COUNT-MIN( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i$ 
    
```

# COUNT-MIN - worked example

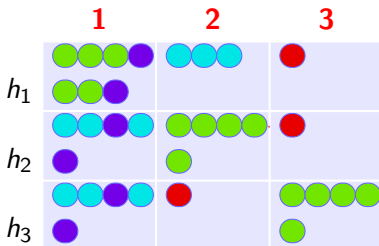


	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

```

COUNT-MIN( $a_i$ )
for each  $j \in [t]$ 
     $C[j, h_j(a_i)] += c_i$ 
    
```

# COUNT-MIN - worked example



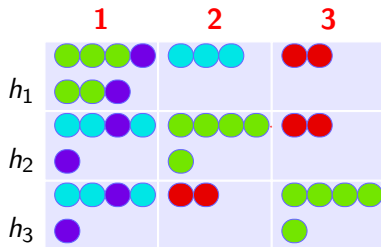
	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

COUNT-MIN( $a_i$ )

for each  $j \in [t]$

$C[j, h_j(a_i)] += c_i$

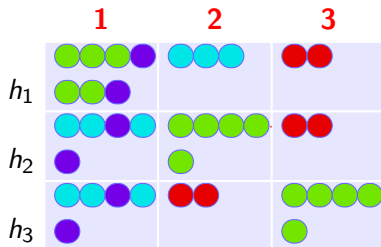
# COUNT-MIN - worked example



	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

return  $\hat{f}_a = \min_{1 \leq i \leq t} C[i, h_i(a)]$

# COUNT-MIN - worked example

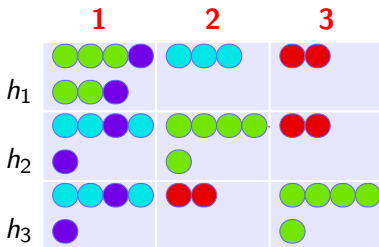


	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

return  $\hat{f}_a = \min_{1 \leq i \leq t} C[i, h_i(a)]$

$$\hat{f}_{\text{green}} = \min(C[1, h_1(\text{green})], C[2, h_2(\text{green})], C[3, h_3(\text{green})]) = \min(7, 5, 5) = 5 \checkmark$$

# COUNT-MIN - worked example



	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

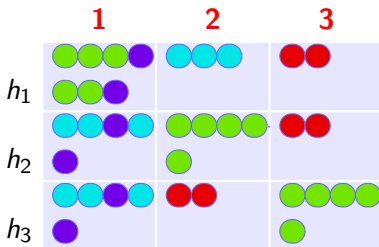
return  $\hat{f}_a = \min_{1 \leq i \leq t} C[i, h_i(a)]$

$$\hat{f}_{\text{green}} = \min(C[1, h_1(\text{green})], C[2, h_2(\text{green})], C[3, h_3(\text{green})]) = \min(7, 5, 5) = 5 \checkmark$$

$$\hat{f}_{\text{cyan}} = \min(C[1, h_1(\text{cyan})], C[2, h_2(\text{cyan})], C[3, h_3(\text{cyan})]) = \min(3, 5, 5) = 3 \checkmark$$



# COUNT-MIN - worked example



	$h_1$	$h_2$	$h_3$
	1	2	3
	2	1	1
	1	1	1
	3	3	2

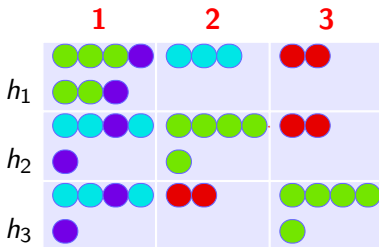
return  $\hat{f}_a = \min_{1 \leq i \leq t} C[i, h_i(a)]$

$$\hat{f}_{\text{green}} = \min(C[1, h_1(\text{green})], C[2, h_2(\text{green})], C[3, h_3(\text{green})]) = \min(7, 5, 5) = 5 \checkmark$$

$$\hat{f}_{\text{cyan}} = \min(C[1, h_1(\text{cyan})], C[2, h_2(\text{cyan})], C[3, h_3(\text{cyan})]) = \min(3, 5, 5) = 3 \checkmark$$

$$\hat{f}_{\text{purple}} = \min(C[1, h_1(\text{purple})], C[2, h_2(\text{purple})], C[3, h_3(\text{purple})]) = \min(7, 5, 5) = 5$$

# COUNT-MIN - worked example



	$h_1$	$h_2$	$h_3$
green	1	2	3
cyan	2	1	1
purple	1	1	1
red	3	3	2

return  $\hat{f}_a = \min_{1 \leq i \leq t} C[i, h_i(a)]$

$$\hat{f}_{\text{green}} = \min(C[1, h_1(\text{green})], C[2, h_2(\text{green})], C[3, h_3(\text{green})]) = \min(7, 5, 5) = 5 \checkmark$$

$$\hat{f}_{\text{cyan}} = \min(C[1, h_1(\text{cyan})], C[2, h_2(\text{cyan})], C[3, h_3(\text{cyan})]) = \min(3, 5, 5) = 3 \checkmark$$

$$\hat{f}_{\text{purple}} = \min(C[1, h_1(\text{purple})], C[2, h_2(\text{purple})], C[3, h_3(\text{purple})]) = \min(7, 5, 5) = 5$$

$$\hat{f}_{\text{red}} = \min(C[1, h_1(\text{red})], C[2, h_2(\text{red})], C[3, h_3(\text{red})]) = \min(2, 2, 2) = 2 \checkmark$$

## COUNT-MIN - Analysis I

For simplicity, consider positive counts of tokens (the cash register model) so that  $\hat{f}_a \geq f_a$  for all tokens  $a$ .

## COUNT-MIN - Analysis I

For simplicity, consider positive counts of tokens (the cash register model) so that  $\hat{f}_a \geq f_a$  for all tokens  $a$ .

Let  $Y_{i,j} = 1$  if  $h_i(j) = h_i(a)$  and 0 otherwise. Note that token  $j$  contributes to  $C[i, h_i(a)]$  iff  $Y_{i,j} = 1$ .

## COUNT-MIN - Analysis I

For simplicity, consider positive counts of tokens (the cash register model) so that  $\hat{f}_a \geq f_a$  for all tokens  $a$ .

Let  $Y_{i,j} = 1$  if  $h_i(j) = h_i(a)$  and 0 otherwise. Note that token  $j$  contributes to  $C[i, h_i(a)]$  iff  $Y_{i,j} = 1$ .

Let r.v.  $X_i$  be the excess count in cell  $C[i, h_i(a)]$ . That is

$$X_i = \sum_{j \in [n] \setminus \{a\}} f_j Y_{i,j}$$

## COUNT-MIN - Analysis I

For simplicity, consider positive counts of tokens (the cash register model) so that  $\hat{f}_a \geq f_a$  for all tokens  $a$ .

Let  $Y_{i,j} = 1$  if  $h_i(j) = h_i(a)$  and 0 otherwise. Note that token  $j$  contributes to  $C[i, h_i(a)]$  iff  $Y_{i,j} = 1$ .

Let r.v.  $X_i$  be the excess count in cell  $C[i, h_i(a)]$ . That is

$$X_i = \sum_{j \in [n] \setminus \{a\}} f_j Y_{i,j}$$

$$\mathbb{E}(X_i) = \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}(Y_{i,j}) = \sum_{j \in [n] \setminus \{a\}} \frac{f_j}{k} = \frac{\|\mathbf{f}\|_1 - f_a}{k} = \frac{\|\mathbf{f}_{-a}\|_1}{k}$$

## COUNT-MIN - Analysis I

For simplicity, consider positive counts of tokens (the cash register model) so that  $\hat{f}_a \geq f_a$  for all tokens  $a$ .

Let  $Y_{i,j} = 1$  if  $h_i(j) = h_i(a)$  and 0 otherwise. Note that token  $j$  contributes to  $C[i, h_i(a)]$  iff  $Y_{i,j} = 1$ .

Let r.v.  $X_i$  be the excess count in cell  $C[i, h_i(a)]$ . That is

$$X_i = \sum_{j \in [n] \setminus \{a\}} f_j Y_{i,j}$$

$$\mathbb{E}(X_i) = \sum_{j \in [n] \setminus \{a\}} f_j \mathbb{E}(Y_{i,j}) = \sum_{j \in [n] \setminus \{a\}} \frac{f_j}{k} = \frac{\|\mathbf{f}\|_1 - f_a}{k} = \frac{\|\mathbf{f}_{-a}\|_1}{k}$$

By Markov's inequality

$$\Pr(X_i \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \frac{\|\mathbf{f}_{-a}\|_1}{k\epsilon \|\mathbf{f}_{-a}\|_1} = \frac{1}{2} \quad \text{set } k = 2/\epsilon$$

## COUNT-MIN - Analysis II

We have a bound for a single counter. Over  $t$  counters the reported excess is the minimum over all  $X_j$ . We can now derive the probability that all the excesses are at least  $\epsilon \|\mathbf{f}_{-a}\|_1$  directly.



## COUNT-MIN - Analysis II

We have a bound for a single counter. Over  $t$  counters the reported excess is the minimum over all  $X_j$ . We can now derive the probability that all the excesses are at least  $\epsilon \|\mathbf{f}_{-a}\|_1$  directly.

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \frac{1}{2^t} = \delta \qquad \text{set } t = \left\lceil \log_2\left(\frac{1}{\delta}\right) \right\rceil$$

## COUNT-MIN - Analysis II

We have a bound for a single counter. Over  $t$  counters the reported excess is the minimum over all  $X_j$ . We can now derive the probability that all the excesses are at least  $\epsilon \|\mathbf{f}_{-a}\|_1$  directly.

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \frac{1}{2^t} = \delta \quad \text{set } t = \left\lceil \log_2 \left( \frac{1}{\delta} \right) \right\rceil$$

$k = 2/\epsilon, t = \lceil \log_2(1/\delta) \rceil$  gives total space in bits

$$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} \cdot (\log m + \log n)\right)$$

## COUNT-MIN - Analysis II

We have a bound for a single counter. Over  $t$  counters the reported excess is the minimum over all  $X_j$ . We can now derive the probability that all the excesses are at least  $\epsilon \|\mathbf{f}_{-a}\|_1$  directly.

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \frac{1}{2^t} = \delta \quad \text{set } t = \left\lceil \log_2 \left( \frac{1}{\delta} \right) \right\rceil$$

$k = 2/\epsilon, t = \lceil \log_2(1/\delta) \rceil$  gives total space in bits

$$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} \cdot (\log m + \log n)\right)$$

The space usage is better than COUNTSKETCH by a factor of  $1/\epsilon$ .

## COUNT-MIN - Analysis II

We have a bound for a single counter. Over  $t$  counters the reported excess is the minimum over all  $X_j$ . We can now derive the probability that all the excesses are at least  $\epsilon \|\mathbf{f}_{-a}\|_1$  directly.

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \frac{1}{2^t} = \delta \quad \text{set } t = \left\lceil \log_2\left(\frac{1}{\delta}\right) \right\rceil$$

$k = 2/\epsilon, t = \lceil \log_2(1/\delta) \rceil$  gives total space in bits

$$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} \cdot (\log m + \log n)\right)$$

The space usage is better than COUNTSKETCH by a factor of  $1/\epsilon$ .

COUNT-MIN's error probability is bounded by  $\epsilon \|\mathbf{f}_{-a}\|_1$  instead of  $\epsilon \|\mathbf{f}_{-a}\|_2$  for COUNTSKETCH.

## COUNT-MIN - Analysis II

We have a bound for a single counter. Over  $t$  counters the reported excess is the minimum over all  $X_j$ . We can now derive the probability that all the excesses are at least  $\epsilon \|\mathbf{f}_{-a}\|_1$  directly.

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \frac{1}{2^t} = \delta \quad \text{set } t = \left\lceil \log_2 \left( \frac{1}{\delta} \right) \right\rceil$$

$k = 2/\epsilon, t = \lceil \log_2(1/\delta) \rceil$  gives total space in bits

$$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} \cdot (\log m + \log n)\right)$$

The space usage is better than COUNTSKETCH by a factor of  $1/\epsilon$ .

COUNT-MIN's error probability is bounded by  $\epsilon \|\mathbf{f}_{-a}\|_1$  instead of  $\epsilon \|\mathbf{f}_{-a}\|_2$  for COUNTSKETCH.

For all vectors  $z \in \mathbb{R}^n$ , we have that  $\|z\|_1 \geq \|z\|_2$ .

## Frequency estimation - space/time summary

We have seen two one-pass sketching algorithms for frequency estimation.

## Frequency estimation - space/time summary

We have seen two one-pass sketching algorithms for frequency estimation.

COUNTSKETCH runs in  $O(t) = O(\log \frac{1}{\delta})$  time per token if  $t = \lceil 1/\delta \rceil$ .

## Frequency estimation - space/time summary

We have seen two one-pass sketching algorithms for frequency estimation.

COUNTSKETCH runs in  $O(t) = O(\log \frac{1}{\delta})$  time per token if  $t = \lceil 1/\delta \rceil$ .

COUNTSKETCH space usage is

$$O(t \log n + tk \log m) = O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right) (\log m + \log n) \text{ bits}$$

bits if  $k = \lceil 3/\epsilon^2 \rceil$ .



## Frequency estimation - space/time summary

We have seen two one-pass sketching algorithms for frequency estimation.

COUNTSKETCH runs in  $O(t) = O(\log \frac{1}{\delta})$  time per token if  $t = \lceil 1/\delta \rceil$ .

COUNTSKETCH space usage is

$$O(t \log n + tk \log m) = O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right) (\log m + \log n) \text{ bits}$$

bits if  $k = \lceil 3/\epsilon^2 \rceil$ .

COUNT-MIN runs in  $O(t) = O(\log \frac{1}{\delta})$  time per token if  $t = \lceil 1/\delta \rceil$ .

## Frequency estimation - space/time summary

We have seen two one-pass sketching algorithms for frequency estimation.

COUNTSKETCH runs in  $O(t) = O(\log \frac{1}{\delta})$  time per token if  $t = \lceil 1/\delta \rceil$ .

COUNTSKETCH space usage is

$$O(t \log n + tk \log m) = O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right) (\log m + \log n) \text{ bits}$$

bits if  $k = \lceil 3/\epsilon^2 \rceil$ .

COUNT-MIN runs in  $O(t) = O(\log \frac{1}{\delta})$  time per token if  $t = \lceil 1/\delta \rceil$ .

COUNT-MIN space usage is

$$O(t \log n + tk \log m) = O\left(\frac{1}{\epsilon} \log\left(\frac{1}{\delta}\right)\right) (\log m + \log n)$$

bits if  $k = \lceil 2/\epsilon \rceil$ . This is a factor of  $1/\epsilon$  improvement.

## Frequency estimation - estimation error summary

COUNTSKETCH: with  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \delta$$

## Frequency estimation - estimation error summary

COUNTSKETCH: with  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \delta$$

COUNT-MIN: with  $k = \lceil 2/\epsilon \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \delta$$

## Frequency estimation - estimation error summary

COUNTSKETCH: with  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \delta$$

COUNT-MIN: with  $k = \lceil 2/\epsilon \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \delta$$

For all vectors  $z \in \mathbb{R}^n$ , we have that  $\|z\|_1 \geq \|z\|_2$  so the estimation error is worse for COUNT-MIN.

## Frequency estimation - estimation error summary

COUNTSKETCH: with  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \delta$$

COUNT-MIN: with  $k = \lceil 2/\epsilon \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \delta$$

For all vectors  $z \in \mathbb{R}^n$ , we have that  $\|z\|_1 \geq \|z\|_2$  so the estimation error is worse for COUNT-MIN.

By setting  $k = 1/\epsilon$ , MISRA-GRIES gives us an estimate

$$f_j - \epsilon \|\mathbf{f}\|_1 \leq \hat{f}_j \leq f_j \text{ for every } j \in [n]$$

## Frequency estimation - estimation error summary

COUNTSKETCH: with  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \delta$$

COUNT-MIN: with  $k = \lceil 2/\epsilon \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \delta$$

For all vectors  $z \in \mathbb{R}^n$ , we have that  $\|z\|_1 \geq \|z\|_2$  so the estimation error is worse for COUNT-MIN.

By setting  $k = 1/\epsilon$ , MISRA-GRIES gives us an estimate

$$f_j - \epsilon \|\mathbf{f}\|_1 \leq \hat{f}_j \leq f_j \text{ for every } j \in [n]$$

MISRA-GRIES gives a lower bound on frequency where COUNT-MIN gives an upper bound.

## Frequency estimation - estimation error summary

COUNTSKETCH: with  $k = \lceil 3/\epsilon^2 \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_2) \leq \delta$$

COUNT-MIN: with  $k = \lceil 2/\epsilon \rceil$  and  $t = \lceil \log_2(1/\delta) \rceil$ ,

$$\Pr(\hat{f}_a - f_a \geq \epsilon \|\mathbf{f}_{-a}\|_1) \leq \delta$$

For all vectors  $z \in \mathbb{R}^n$ , we have that  $\|z\|_1 \geq \|z\|_2$  so the estimation error is worse for COUNT-MIN.

By setting  $k = 1/\epsilon$ , MISRA-GRIES gives us an estimate

$$f_j - \epsilon \|\mathbf{f}\|_1 \leq \hat{f}_j \leq f_j \text{ for every } j \in [n]$$

MISRA-GRIES gives a lower bound on frequency where COUNT-MIN gives an upper bound.

MISRA-GRIES uses  $O((1/\epsilon)(\log m + \log n))$  bits but does not work in the turnstile model (with deletions).