# CountMin sketch

Ben Langmead

JOHNS HOPKINS

WHITING SCHOOL
*of* ENGINEERING

## Department of Computer Science

# Counting

Input is a "stream" of items $\{a_1, a_2, \ldots, a_m\}$, each from universe of size $n$.

Number of times a value $x$ appears is its "count" or "frequency" $f_x$

Stream of zip-code digits: 2, 1, 2, 1, 8, 2, 6, 8, 2

$$m = 9$$

$$n = |\{0,1,\ldots,9\}| = 10$$

$$f_1 = 2 \qquad f_2 = 4$$

# Aside on notation

Defining variables like $n, m, N, M,$ and having to specify "distinct" versus not, can get tiresome

An alternative is to pick a variable for the input data stream, say $\mathbf{a}$

Then use double bars to express "moments"

$$\| \mathbf{a} \|_1 = \sum_{x \in distinct(\mathbf{a})} f_x \qquad \| \mathbf{a} \|_0 = \sum_{x \in distinct(\mathbf{a})} (f_x)^0$$

\# items in stream            # *distinct* items in stream

# Aside on notation

We can also consider higher moments like $2$:
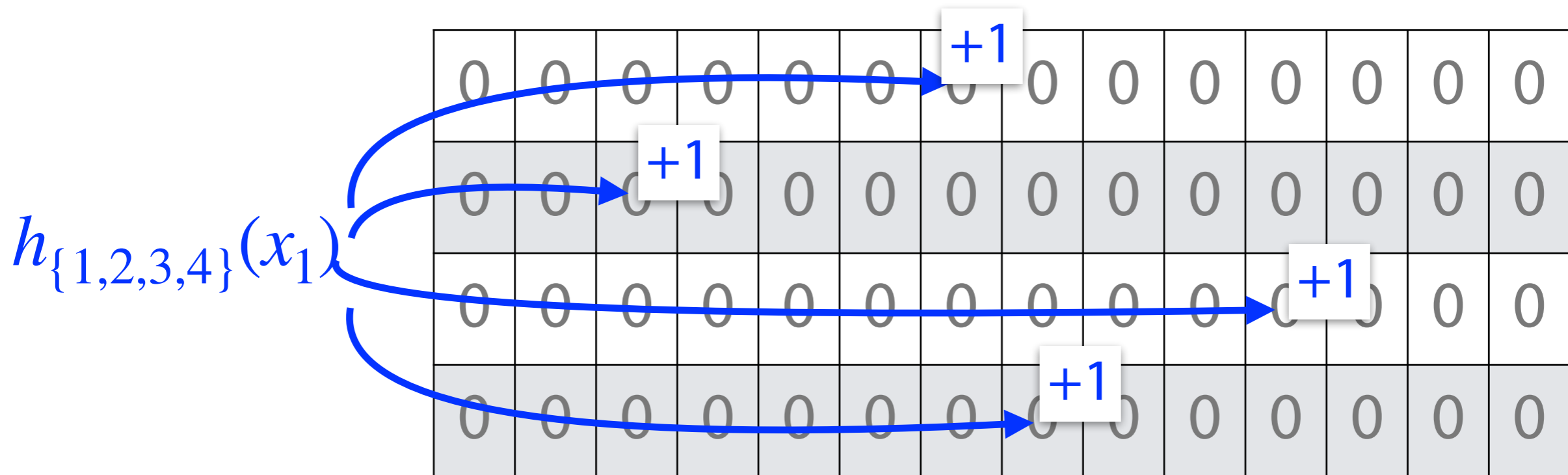
$$\| \mathbf{a} \|_2 = \sum_{x \in distinct(\mathbf{a})} \left( f_x \right)^2$$

Or, more generally, $k$:

$$\| \mathbf{a} \|_k = \sum_{x \in distinct(\mathbf{a})} \left( f_x \right)^k$$

Today we're concerned with $f_x$, not its powers

# CountMin

Matrix of counters, all initially 0

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# CountMin

Insert:

$h_{\{1,2,3,4\}}(x_1)$

# CountMin

| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

# CountMin

Insert:

$h_{\{1,2,3,4\}}(x_2)$

| 0 | 0 | 0 +1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 +1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 +1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 +1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

# CountMin

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

# CountMin

Insert:

$h_{\{1,2,3,4\}}(x_3)$

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 +1 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 +1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 +1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 +1 | 0 | 0 | 0 | 0 | 0 | 0 |

# CountMin

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |

# CountMin

Point query:

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |

$h_{\{1,2,3,4\}}(q_1)$

What should the estimate $\tilde{f}_x$ be?     0

# CountMin

Point query:

$h_{\{1,2,3,4\}}(q_1)$

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 |   |   |   |   |   |   | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 |   |   |   |   |   | 3 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |

**Definite collisions** in these two rows

How much is due to collisions?

# CountMin

Point query:

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |

$h_{\{1,2,3,4\}}(q_2)$

What should the estimate $\tilde{f}_x$ be?    1

# CountMin

Point query:

$h_{\{1,2,3,4\}}(q_2)$

| 0 | 0 | 1 | | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 | | 0 | 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | | 0 | 0 | 0 | 3 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 0 | |

Could be collisions?

**Yes**

**Definite collision** in this row

How much is due to collisions?

# CountMin



Query for item $x$ returns **_minimum_** of the selected elements; call this estimate $\tilde{f}_x$

Collisions can make us overestimate $f_x$, but not underestimate; i.e. $$\tilde{f}_x \geq f_x$$

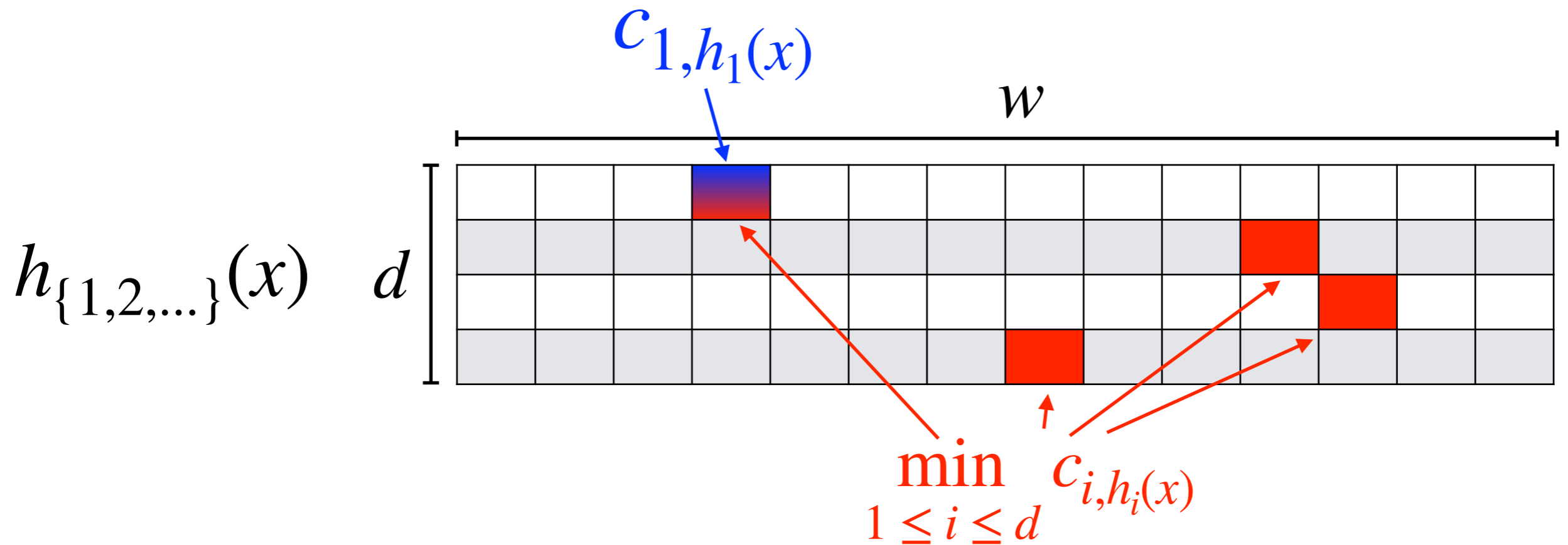Can we argue $\tilde{f}_x$ is **_probably not too far_** from $f_x$?

# CountMin



We use functions $h_1, h_2, \ldots, h_d$ drawn from family $H$, each ranging over $\{1, 2, \ldots, w\}$

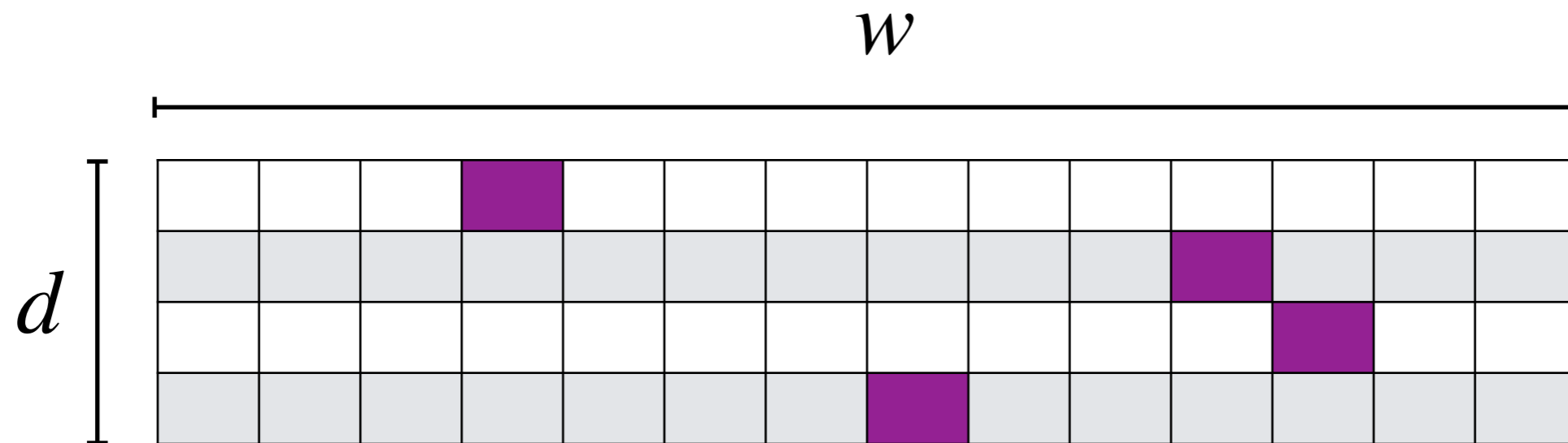Let $c_{i,j}$ be the number of items $x$ such that $h_i(x) = j$

# CountMin



What's the relationship between: $f_x$ $\textcolor{blue}{c_{1,h_1(x)}}$ $\textcolor{red}{\min_{1 \le i \le d} c_{i,h_i(x)}}$

$$f_x \le \textcolor{red}{\min_{1 \le i \le d} c_{i,h_i(x)}} \le \textcolor{blue}{c_{1,h_1(x)}}$$

# CountMin

$$w$$



$$d$$

Recall $m$ = # of items in stream

Claim: if $w = 2/\epsilon$ and $d = \log_2 \delta^{-1}$, then

$$\Pr\left(\tilde{f}_x \leq f_x + \epsilon m\right) \geq 1 - \delta$$

$\tilde{f}_x$ is **probably** **not too far** from $f_x$

# CountMin

Pick item $x$ and define r.v.s $\{Z_1, Z_2, \ldots, Z_d\}$ such that $Z_i = c_{i,h_i(x)} - f_x$

$Z_i$ is the amount we over-counted in row $i$ due to collisions

# CountMin

For $i \in \{1, 2, \ldots, d\}, y \in \{\text{distinct items}\} \setminus \{x\}$

$$X_{i,y} = \begin{cases} 1 & \text{if } h_i(y) = h_i(x) \\ 0 & \text{otherwise} \end{cases}$$

Recall: $Z_i$ is the amount we over-counted in row $i$ due to collisions

$$Z_i = \sum_{y \neq x} \left( f_y \cdot X_{i,y} \right)$$

# CountMin

$$\mathbf{E}[Z_i] = \mathbf{E}\left[\sum_{x \neq y} f_y \cdot X_{i,y}\right]$$

$$= \sum_{x \neq y} f_y \cdot \mathbf{E}\left[X_{i,y}\right] \qquad \text{Linearity of expectation}$$

$$= \sum_{x \neq y} f_y \cdot \Pr\left(h_i(y) = h_i(x)\right) \qquad \text{Expectation of indicator}$$

What would we like to use next?   **2-universality**

# CountMin

Further assume that family $H$ from which $h_i$'s were drawn is 2-universal

$$\sum_{x \neq y} f_y \cdot \text{Pr}\left(h_i(y) = h_i(x)\right) \leq \sum_{x \neq y} f_y \cdot \frac{1}{w} \quad \text{2-universality}$$

$$\leq \frac{m}{w}$$

Expected per-row excess $\mathbf{E}[Z_i]$ is at most $m/w$

# CountMin

$Z_i$ is a non-negative r.v., so:

$$\Pr\left(Z_i \geq a\right) \leq \frac{\mathbf{E}[Z_i]}{a}$$

Markov inequality

Let $b = \dfrac{a}{\mathbf{E}[Z_i]}$

$$\Pr\left(Z_i \geq b \cdot \mathbf{E}[Z_i]\right) \leq \frac{1}{b}$$

# CountMin

$$\Pr\left(Z_i \geq b \cdot \mathbf{E}[Z_i]\right) \leq \frac{1}{b}$$

Combine with $\mathbf{E}[Z_i] \leq \dfrac{m}{w}$ :

$$\Pr\left(Z_i \geq \frac{bm}{w}\right) \leq \Pr\left(Z_i \geq b \cdot \mathbf{E}[Z_i]\right) \leq \frac{1}{b}$$

Continue with these

# CountMin

$$\Pr\left(Z_i \geq \frac{bm}{w}\right) \leq \frac{1}{b}$$

Let $b = w\epsilon$ (# columns times error tolerance):

$$\Pr\left(Z_i \geq \epsilon m\right) \leq \frac{1}{w\epsilon}$$

Let $w = 2/\epsilon$ (# columns from our Claim):

$$\Pr\left(Z_i \geq \epsilon m\right) \leq \frac{1}{2}$$

# CountMin

When $w = 2/\epsilon$, probability that "bad thing" happens is at most 1/2

$$\Pr\left(Z_i \geq \epsilon m\right) = \Pr\left(f_x + Z_i \geq f_x + \epsilon m\right) \leq \frac{1}{2}$$

We want an upper bound of $\delta$, $\delta$ being small

So: Repeat (across rows) and take minimum

# CountMin

$$\Pr\left(Z_i \geq \epsilon m\right) \leq \frac{1}{2}$$

$$\Pr\left(\forall_{1 \leq i \leq d} \; Z_i \geq \epsilon m\right) \leq \left(\frac{1}{2}\right)^d$$

<span style="color:blue">Independence of uniform & independently chosen hashes</span>

Recall we set $d = \log_2 \delta^{-1}$

$$\left(\frac{1}{2}\right)^d = 2^{-\log_2 \delta^{-1}} = 2^{\log_2 \delta} = \delta$$

# CountMin

$$\Pr\left(\forall_{1 \le i \le d} \ Z_i \ge \epsilon m\right) \le \delta$$

Probability of the bad thing is $\le \delta$

$\downarrow$ complement

$$\Pr\left(\exists_{1 \le i \le d} \ Z_i < \epsilon m\right) \ge 1 - \delta$$

Prob. of good thing is $\ge 1 - \delta$

To get the good thing, take the minimum

# CountMin
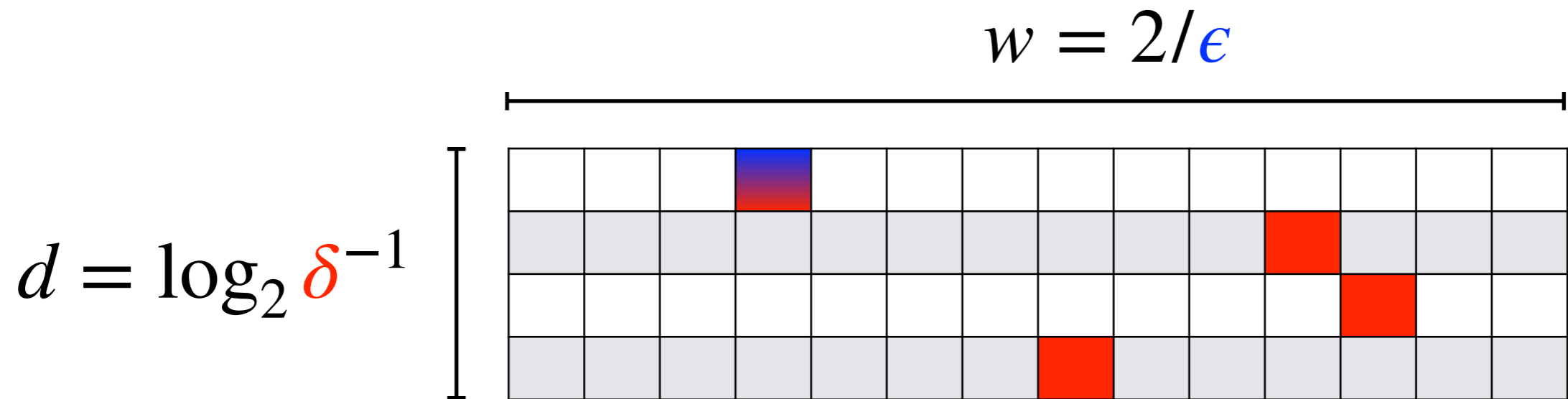
Claim is proved:

$$\tilde{f}_x = \min(c_{1,h_1(x)}, c_{2,h_2(x)}, \ldots, c_{d,h_d(x)})$$

$$= \min(f_x + Z_1, f_x + Z_2, \ldots, f_x + Z_d) \leq f_x + \epsilon m$$

With probability $1 - \delta$ ✅

$\tilde{f}_x$ is **probably not too far** from $f_x$

# CountMin

$$w = 2/\textcolor{blue}{\epsilon}$$

$$d = \log_2 \textcolor{red}{\delta}^{-1}$$

To achieve this, sketch must contain
$O(\epsilon^{-1} \log \delta^{-1})$ counters

# CountMin

| $\epsilon$ | $\delta$ | $\lceil (2/\epsilon) \rceil \cdot \lceil \log_2 \delta^{-1} \rceil$ |
|---:|---:|---:|
| 10% | 0.1 | 80 |
| 1% | 0.01 | 1,400 |
| 0.1% | 0.001 | 20,000 |
| 0.0001% | 0.01 | 1,400,000 |

Remember that $\epsilon$ multiplies $m$, and a counter requires many (maybe 32 or 64) bits

# CountMin

A common use of CountMin is to find **heavy hitters**

Items with frequency over a threshold

$h_{\{1,2,3,4\}}(x)$

| 10 | 5 | 17 | 17 | 17 | 1 | 8 | 20 | 18 | 14 | 8 | 5 | 20 | 13 |
| 15 | 4 | 4 | 19 | 20 | 0 | 8 | 17 | 15 | 19 | 4 | 20 | 16 | 12 |
| 18 | 14 | 9 | 10 | 10 | 15 | 10 | 9 | 16 | 4 | 10 | 18 | 20 | 10 |
| 13 | 3 | 6 | 18 | 8 | 19 | 1 | 15 | 11 | 1 | 8 | 8 | 18 | 3 |

+1 (on 20)
+1 (on 19)
+1 (on 10)
+1 (on 15)

While adding items, maintain data structure containing items with point query result over the threshold

# CountMin

$$\left[\quad M \quad\right]\left[\begin{array}{c} \\ v \\ \\ \end{array}\right] = \left[\; Mv \;\right] \longrightarrow \text{answer}$$

Linear transformation
interpretation

$v$ is input data
$M$ builds the sketch

# CountMin

Result of applying $h_1$ to $x_1, x_2, \ldots, x_8$

Row 1 of CountMin

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} f_{x_1} \\ f_{x_2} \\ f_{x_3} \\ f_{x_4} \\ f_{x_5} \\ f_{x_6} \\ f_{x_7} \\ f_{x_8} \end{bmatrix} = \begin{bmatrix} f_{x_3} + f_{x_4} \\ f_{x_1} + f_{x_5} + f_{x_7} \\ f_{x_2} + f_{x_6} + f_{x_8} \end{bmatrix}$$

Adapted from Andrew McGregor:
https://people.cs.umass.edu/~mcgregor/stocworkshop/mcgregor.pdf