

Count Sketches and Bloom Filters

Recap from Last Time

Distribution Property:

Each element should have an equal probability of being placed in each slot.

For any $x \in \mathcal{U}$ and random $h \in \mathcal{H}$, the value of $h(x)$ is uniform over $[m]$.

Independence Property:

Where one element is placed shouldn't impact where a second goes.

For any distinct $x, y \in \mathcal{U}$ and random $h \in \mathcal{H}$, $h(x)$ and $h(y)$ are independent random variables.

A family of hash functions \mathcal{H} is called ***2-independent*** (or ***pairwise independent***) if it satisfies the distribution and independence properties.

Suppose there are two tunable values

$$\varepsilon \in (0, 1]$$

$$\delta \in (0, 1]$$

where ε represents *accuracy* and δ represents *confidence*.

Goal: Make an estimator \hat{A} for some quantity A where

With probability at least $1 - \delta$,
 $|\hat{A} - A| \leq \varepsilon \cdot \text{size}(\text{input})$

Probably
Approximately Correct

for some measure of the size of the input.

What does it mean for an approximation to be “good”?

How to Build an Estimator

	<i>Count-Min Sketch</i>
<i>Step One:</i> Build a Simple Estimator	Hash items to counters; add +1 when item seen.
<i>Step Two:</i> Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.
<i>Step Three:</i> Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.
<i>Step Four:</i> Replicate to Boost Confidence	Take min; only fails if all estimates are bad.

New Stuff!

The Count Sketch



Frequency Estimation

- **Recall:** A frequency estimator is a data structure that supports
 - **increment**(x), which increments the number of times that we've seen x , and
 - **estimate**(x), which returns an estimate of how many times we've seen x .
- **Notation:** Assume that the elements we're processing are x_1, \dots, x_n , and that the true frequency of element x_i is a_i .
- Remember that the frequencies are not random variables – we're assuming that they're not under our control. Any randomness comes from hash functions.


How to Build an Estimator



	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

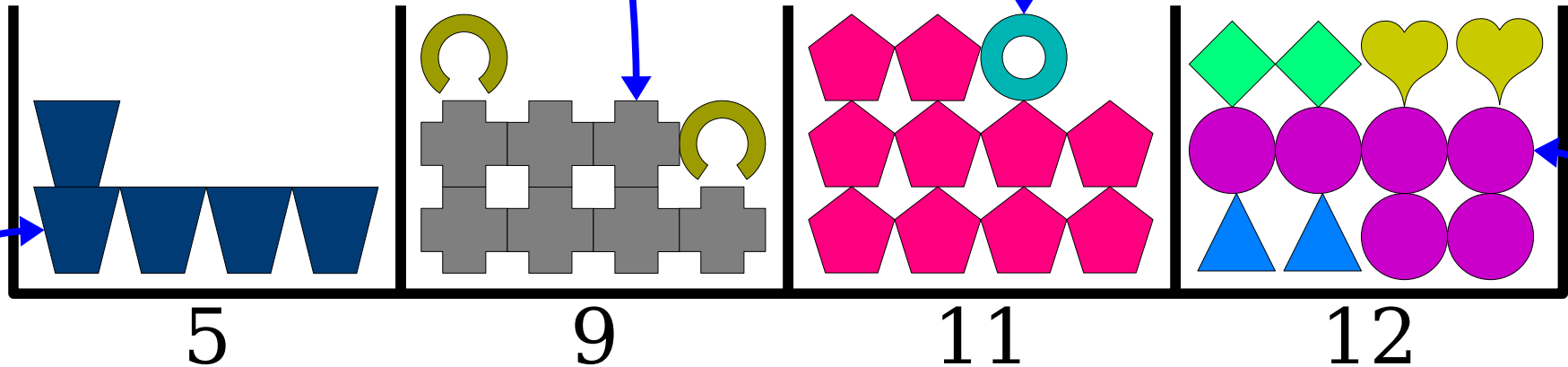
How to Build an Estimator


	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	


Revisiting Count-Min

We have a reasonable estimate for , since it collides with an uncommon item.


No matter what we do, we're not going to get a good estimate for  because it collides with a very frequent item ().





We have a good estimate for , since nothing collides with it.


Our estimate for  is way off because of lots of small collisions.


Revisiting Count-Min

We have a reasonable estimate for , since it collides with an uncommon item.

No matter what we do, we're not going to get a good estimate for  because it collides with a very frequent item ().

Question: Can we mitigate the impact of collisions with lots of infrequent elements?

We have a good estimate for , since nothing collides with it.

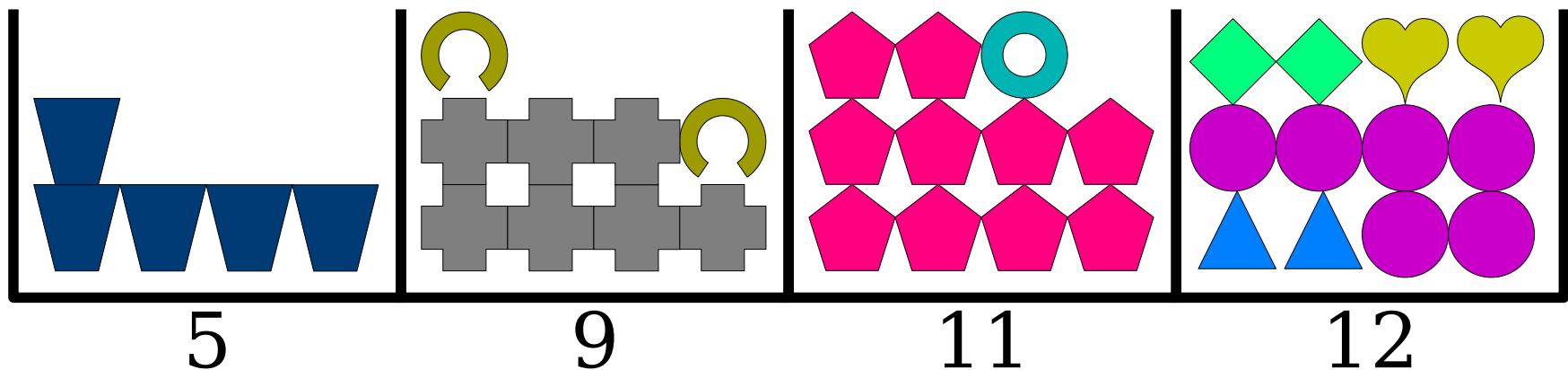
Our estimate for  is way off because of lots of small collisions.

5

12

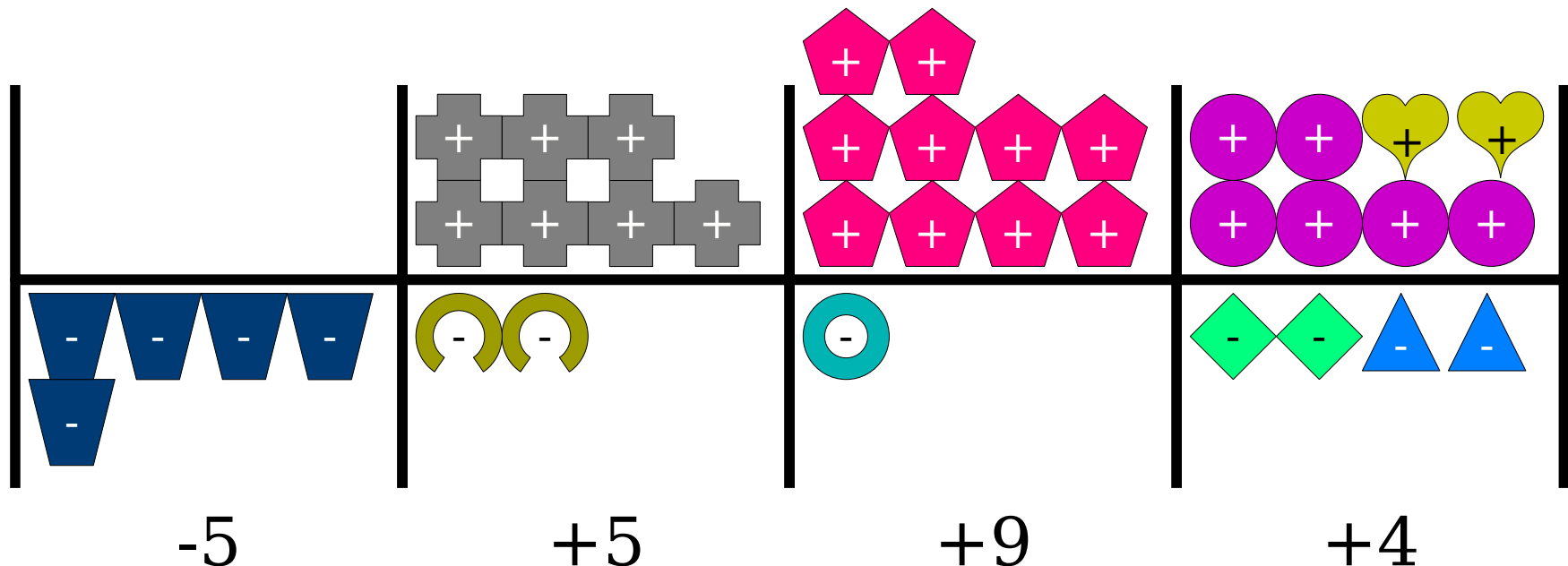
The Setup

- As before, create an array of counters and assign each item a counter.
- **Key New Step:** For each item x , assign x either $+1$ or -1 .



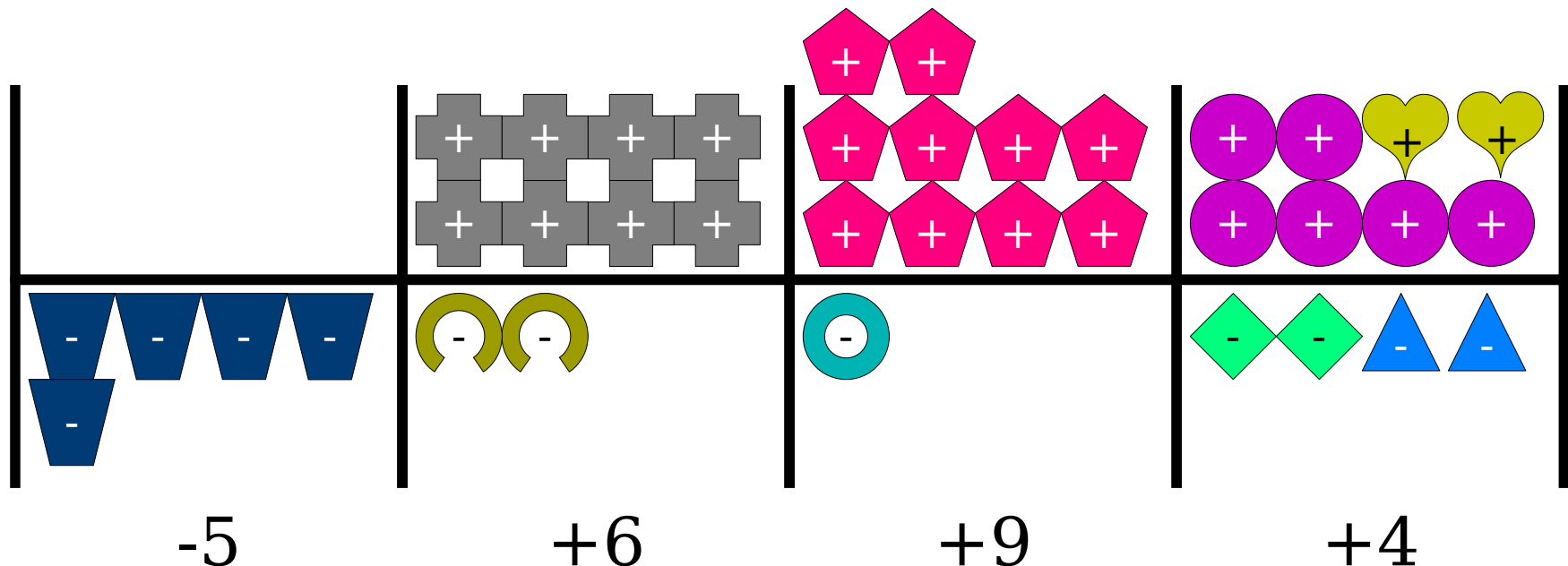
The Setup

- As before, create an array of counters and assign each item a counter.
- **Key New Step:** For each item x , assign x either $+1$ or -1 .
 - To **increment**(x), go to **count**[$h(x)$] and add ± 1 as appropriate.
 - To **estimate**(x), return **count**[$h(x)$], multiplied by ± 1 as appropriate.



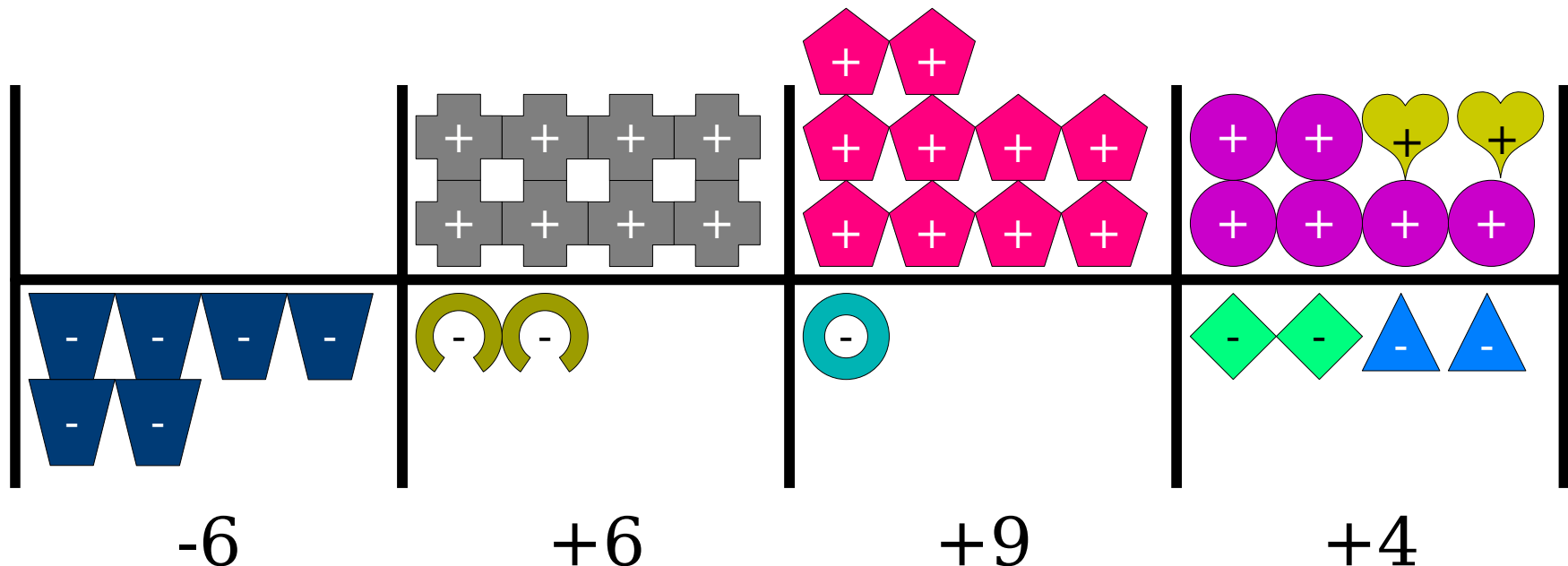
The Setup

- As before, create an array of counters and assign each item a counter.
- **Key New Step:** For each item x , assign x either $+1$ or -1 .
 - To **increment**(x), go to **count**[$h(x)$] and add ± 1 as appropriate.
 - To **estimate**(x), return **count**[$h(x)$], multiplied by ± 1 as appropriate.



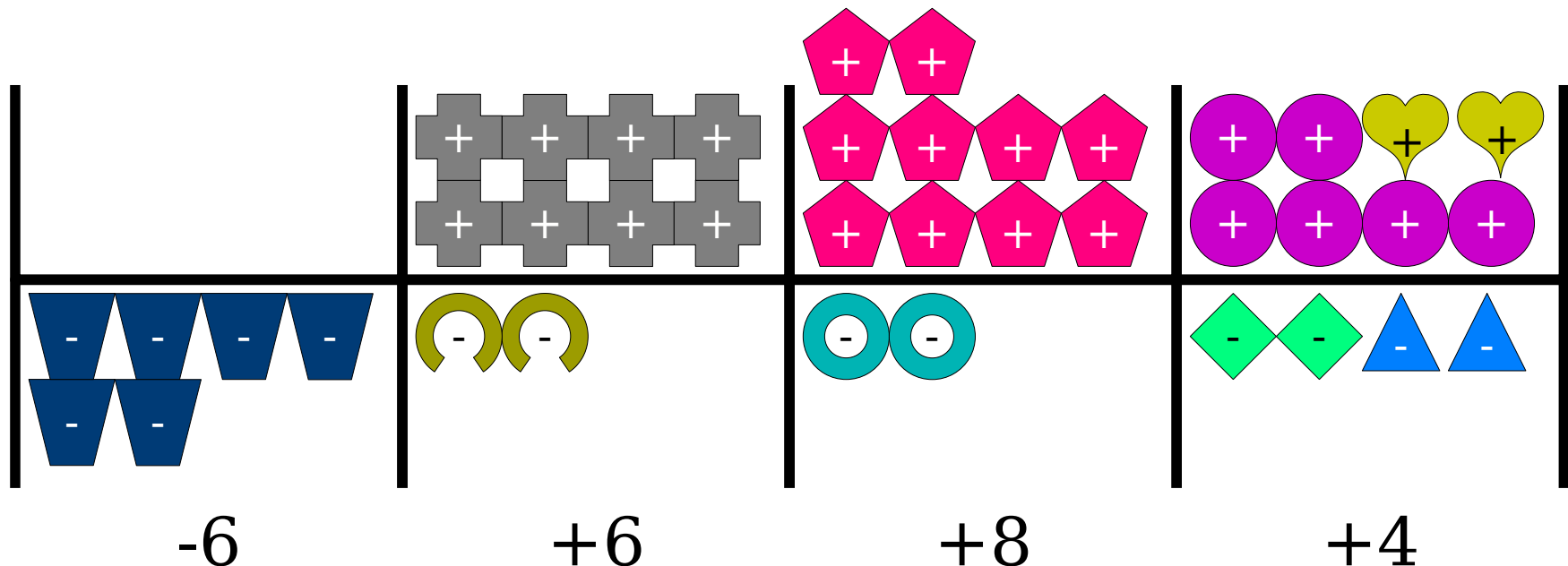
The Setup

- As before, create an array of counters and assign each item a counter.
- **Key New Step:** For each item x , assign x either $+1$ or -1 .
 - To **increment**(x), go to **count**[$h(x)$] and add ± 1 as appropriate.
 - To **estimate**(x), return **count**[$h(x)$], multiplied by ± 1 as appropriate.





The Setup

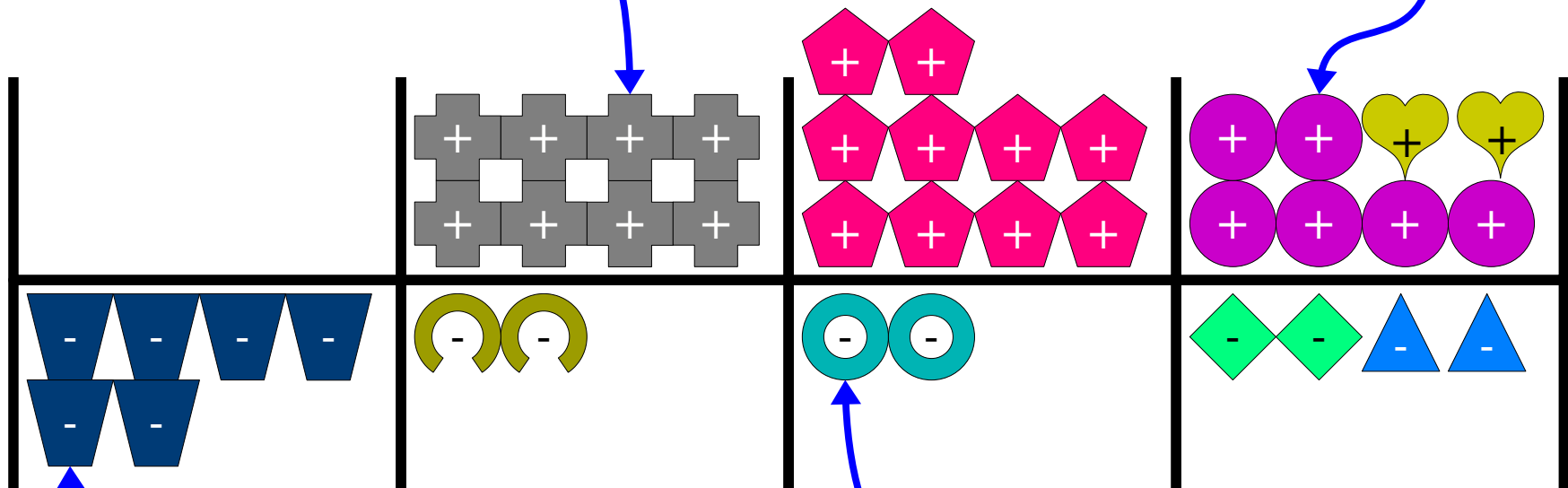
- As before, create an array of counters and assign each item a counter.
- **Key New Step:** For each item x , assign x either $+1$ or -1 .
 - To **increment**(x), go to **count**[$h(x)$] and add ± 1 as appropriate.
 - To **estimate**(x), return **count**[$h(x)$], multiplied by ± 1 as appropriate.






The Setup

We have a reasonable estimate for , since it collides with an uncommon item.

We have a reasonable estimate for  because the other collisions mostly offset.

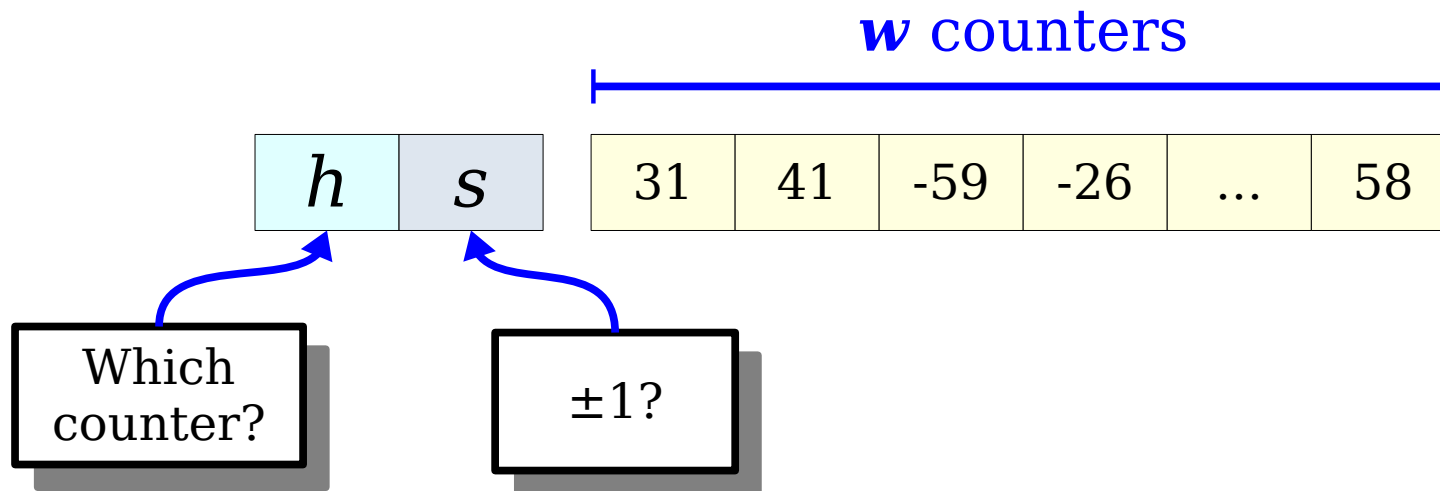


We have a good estimate for , since nothing collides with it.

No matter what we do, we're not going to get a good estimate for  because it collides with a very frequent item ().

Formalizing This

- Maintain an array of counters of length w .
- Pick $h \in \mathcal{H}$ chosen uniformly at random from a 2-independent family of hash functions from \mathcal{U} to w .
- Pick $s \in \mathcal{U}$ uniformly randomly and independently of h from a 2-independent family from \mathcal{U} to $\{-1, +1\}$.
- *increment*(x): $\text{count}[h(x)] += s(x)$.
- *estimate*(x), return $s(x) \cdot \text{count}[h(x)]$.



How to Build an Estimator

	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

How to Build an Estimator

	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
<i>Step One:</i> Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	
<i>Step Three:</i> Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	
<i>Step Four:</i> Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

Formalizing the Intuition

- Define $\hat{\mathbf{a}}_i$ to be our estimate of \mathbf{a}_i .
- As before, $\hat{\mathbf{a}}_i$ will depend on how the other elements are distributed. Unlike before, it now also depends on signs given to the elements by s .
- Specifically, for each other x_j that collides with x_i , the estimate $\hat{\mathbf{a}}_i$ includes an error term of

$$s(x_i) \cdot s(x_j) \cdot \mathbf{a}_j$$

- Why?

Formalizing the Intuition

- Define $\hat{\mathbf{a}}_i$ to be our estimate of \mathbf{a}_i .
- As before, $\hat{\mathbf{a}}_i$ will depend on how the other elements are distributed. Unlike before, it now also depends on signs given to the elements by s .
- Specifically, for each other x_j that collides with x_i , the estimate $\hat{\mathbf{a}}_i$ includes an error term of

$$s(x_i) \cdot s(x_j) \cdot \mathbf{a}_j$$

- Why?
 - The counter for x_i will have $s(x_j) \mathbf{a}_j$ added in.
 - We multiply the counter by $s(x_i)$ before returning it.

Formalizing the Intuition

- Define $\hat{\mathbf{a}}_i$ to be our estimate of \mathbf{a}_i .
- As before, $\hat{\mathbf{a}}_i$ will depend on how the other elements are distributed. Unlike before, it now also depends on signs given to the elements by s .
- Specifically, for each other x_j that collides with x_i , the estimate $\hat{\mathbf{a}}_i$ includes an error term of

$$s(x_i) \cdot s(x_j) \cdot \mathbf{a}_j$$

- Why?
 - If $s(x_i)$ and $s(x_j)$ point in the same direction, the terms add to the total.
 - If $s(x_i)$ and $s(x_j)$ point in different directions, the terms subtract from the total.

Formalizing the Intuition

- In our quest to learn more about $\hat{\mathbf{a}}_i$, let's have X_j be a random variable indicating whether \mathbf{x}_i and \mathbf{x}_j collided with one another:

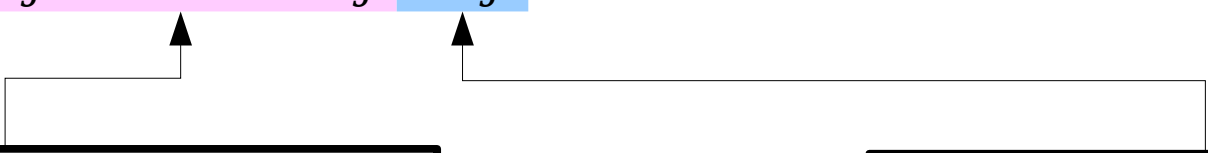
$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

Formalizing the Intuition

- In our quest to learn more about $\hat{\mathbf{a}}_i$, let's have X_j be a random variable indicating whether \mathbf{x}_i and \mathbf{x}_j collided with one another:

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

- We can then express $\hat{\mathbf{a}}_i$ in terms of the signed contributions from the items \mathbf{x}_i collides with:

$$\hat{\mathbf{a}}_i = \sum_j \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j$$


This is how much the collision impacts our estimate.

We only care about items we collided with.

Formalizing the Intuition

- In our quest to learn more about $\hat{\mathbf{a}}_i$, let's have X_j be a random variable indicating whether \mathbf{x}_i and \mathbf{x}_j collided with one another:

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

- We can then express $\hat{\mathbf{a}}_i$ in terms of the signed contributions from the items \mathbf{x}_i collides with:

$$\hat{\mathbf{a}}_i = \sum_j \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j = \mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j$$

This is how much the collision impacts our estimate.

We only care about items we collided with.

$$\mathbb{E}[\hat{\boldsymbol{a}}_i] = \mathbb{E}[\boldsymbol{a}_i + \sum_{j \neq i} \boldsymbol{a}_j s(\boldsymbol{x}_i) s(\boldsymbol{x}_j) X_j]$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
\end{aligned}$$

Hey, it's
linearity of
expectation!

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
\end{aligned}$$

Remember that
 \mathbf{a}_i and the like
 aren't random
 variables.

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
\end{aligned}$$

We chose the hash functions h and s independently of one another.

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

We chose the hash functions h and s independently of one another.

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

Since s is drawn from a 2-independent family of hash functions, we know $s(\mathbf{x}_i)$ and $s(\mathbf{x}_j)$ are independent random variables.

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

$$\mathbb{E}[s(\mathbf{x}_i)] =$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

$$\mathbb{E}[s(\mathbf{x}_i)] =$$

s is drawn from a 2-independent family of hash functions.

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

$$\mathbb{E}[s(\mathbf{x}_i)] =$$

s is drawn from a 2-independent family of hash functions.

$s(\mathbf{x}_i)$ is uniform over $\{-1, +1\}$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

$$\mathbb{E}[s(\mathbf{x}_i)] =$$

s is drawn from a 2-independent family of hash functions.

$s(\mathbf{x}_i)$ is uniform over $\{-1, +1\}$

$$\Pr[s(\mathbf{x}_i) = -1] = 1/2 \quad \Pr[s(\mathbf{x}_i) = +1] = 1/2$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

$$\mathbb{E}[s(\mathbf{x}_i)] = \frac{1}{2} \cdot (-1) + \frac{1}{2} \cdot (+1)$$

s is drawn from a 2-independent family of hash functions.

$s(\mathbf{x}_i)$ is uniform over $\{-1, +1\}$

$$\Pr[s(\mathbf{x}_i) = -1] = \frac{1}{2} \quad \Pr[s(\mathbf{x}_i) = +1] = \frac{1}{2}$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j]
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}[s(\mathbf{x}_i)] &= \frac{1}{2} \cdot (-1) + \frac{1}{2} \cdot (+1) \\
&= 0
\end{aligned}$$

s is drawn from a 2-independent family of hash functions.

$s(\mathbf{x}_i)$ is uniform over $\{-1, +1\}$

$$\Pr[s(\mathbf{x}_i) = -1] = \frac{1}{2} \quad \Pr[s(\mathbf{x}_i) = +1] = \frac{1}{2}$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} 0
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}[s(\mathbf{x}_i)] &= \frac{1}{2} \cdot (-1) + \frac{1}{2} \cdot (+1) \\
&= 0
\end{aligned}$$

s is drawn from a 2-independent family of hash functions.

$s(\mathbf{x}_i)$ is uniform over $\{-1, +1\}$

$$\Pr[s(\mathbf{x}_i) = -1] = \frac{1}{2} \quad \Pr[s(\mathbf{x}_i) = +1] = \frac{1}{2}$$

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i) s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(\mathbf{x}_i)] \mathbb{E}[s(\mathbf{x}_j)] \mathbb{E}[\mathbf{a}_j X_j] \\
&= \mathbf{a}_i + \sum_{j \neq i} 0 \\
&= \mathbf{a}_i
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}[s(\mathbf{x}_i)] &= \frac{1}{2} \cdot (-1) + \frac{1}{2} \cdot (+1) \\
&= 0
\end{aligned}$$

s is drawn from a 2-independent family of hash functions.

$s(\mathbf{x}_i)$ is uniform over $\{-1, +1\}$

$$\Pr[s(\mathbf{x}_i) = -1] = \frac{1}{2} \quad \Pr[s(\mathbf{x}_i) = +1] = \frac{1}{2}$$

How to Build an Estimator

	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	2-independence breaks up products; ± 1 variables have zero expected value.
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

How to Build an Estimator

	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	2-independence breaks up products; ± 1 variables have zero expected value.
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

A Hitch

- In the count-min sketch, we used Markov's inequality to bound the probability that we get a bad estimate.
- This worked because we had a ***one-sided error***: the distance $\hat{\mathbf{a}}_i - \mathbf{a}_i$ from the true answer was nonnegative.
- With the count sketch, we have a ***two-sided error***: $\hat{\mathbf{a}}_i - \mathbf{a}_i$ can be negative in the count sketch because collisions can *decrease* the estimate $\hat{\mathbf{a}}_i$ below the true value \mathbf{a}_i .
- We'll need to use a different technique to bound the error.

Chebyshev to the Rescue

- ***Chebyshev's inequality*** states that for any random variable X with finite variance, given any $c > 0$, we have

$$\Pr[|X - E[X]| \geq c] \leq \frac{\text{Var}[X]}{c^2}.$$

- If we can get the variance of $\hat{\mathbf{a}}_i$, we can bound the probability that we get a bad estimate with our data structure.

$$\text{Var}[\hat{\boldsymbol{a}}_i] = \text{Var}[\boldsymbol{a}_i + \sum_{j \neq i} \boldsymbol{a}_j s(\boldsymbol{x}_i) s(\boldsymbol{x}_j) X_j]$$

$$\text{Var}[\hat{\mathbf{a}}_i] = \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]$$

$$\text{Var}[a + X] = \text{Var}[X]$$

$$\begin{aligned}
 \text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
 &= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
 \end{aligned}$$

$$\text{Var}[a + X] = \text{Var}[X]$$

$$\begin{aligned}\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\ &= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]\end{aligned}$$

In general, Var is *not* a linear operator.

However, if the terms in the sum are ***pairwise uncorrelated***, then Var is linear.

Lemma: The terms in this sum are uncorrelated. (*Prove this!*)

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
\end{aligned}$$

In general, Var is *not* a linear operator.

However, if the terms in the sum are ***pairwise uncorrelated***, then Var is linear.

Lemma: The terms in this sum are uncorrelated. (*Prove this!*)

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
\end{aligned}$$



The “Sum-o’-Var”
Samovar!

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]
\end{aligned}$$

$$\begin{aligned}
\text{Var}[Z] &= \text{E}[Z^2] - \text{E}[Z]^2 \\
&\leq \text{E}[Z^2]
\end{aligned}$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2]
\end{aligned}$$

$$\begin{aligned}
\text{Var}[Z] &= \text{E}[Z^2] - \text{E}[Z]^2 \\
&\leq \text{E}[Z^2]
\end{aligned}$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2]
\end{aligned}$$

$$s(\mathbf{x}) = \pm 1,$$

so

$$s(\mathbf{x})^2 = 1$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2]
\end{aligned}$$

$$s(\mathbf{x}) = \pm 1,$$

so

$$s(\mathbf{x})^2 = 1$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2]
\end{aligned}$$

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2]
\end{aligned}$$

$$X_j^2 = \begin{cases} 1^2 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0^2 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2]
\end{aligned}$$

Useful Fact: If X is an indicator, then $X^2 = X$.

$$X_j^2 = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j]
\end{aligned}$$

Useful Fact: If X is an indicator, then $X^2 = X$.

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j]
\end{aligned}$$

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\text{Var}[\hat{\mathbf{a}}_i] = \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]$$

$$= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]$$

$$= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j]$$

$$\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2]$$

$$= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2]$$

$$= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2]$$

$$= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j]$$

$$= \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2$$

$$X_j = \begin{cases} 1 & \text{if } h(\mathbf{x}_i) = h(\mathbf{x}_j) \\ 0 & \text{if } h(\mathbf{x}_i) \neq h(\mathbf{x}_j) \end{cases}$$

$$\begin{aligned}
\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \text{Var}[\sum_{j \neq i} \mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&= \sum_{j \neq i} \text{Var}[\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j] \\
&\leq \sum_{j \neq i} \text{E}[(\mathbf{a}_j s(\mathbf{x}_i) s(\mathbf{x}_j) X_j)^2] \\
&= \sum_{j \neq i} \text{E}[\mathbf{a}_j^2 s(\mathbf{x}_i)^2 s(\mathbf{x}_j)^2 X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j^2] \\
&= \sum_{j \neq i} \mathbf{a}_j^2 \text{E}[X_j] \\
&= \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2
\end{aligned}$$

I know this might look really dense, but many of these substeps end up being really useful techniques. These ideas generalize, I promise.

Think of $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots]$ as a vector.

What does the following quantity represent?

$$\sum_j \mathbf{a}_j^2$$

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2$$

Think of $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots]$ as a vector.

What does the following quantity represent?

$$\sum_j \mathbf{a}_j^2$$

This is the square of the magnitude of the vector!

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2$$

Think of $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots]$ as a vector.

What does the following quantity represent?

$$\sum_j \mathbf{a}_j^2$$

This is the square of the magnitude of the vector!

The magnitude of a vector is called its **L_2 norm** and is denoted $\|\mathbf{a}\|_2$.

$$\|\mathbf{a}\|_2 = \sqrt{\sum_j \mathbf{a}_j^2}$$

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2$$

Think of $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots]$ as a vector.

What does the following quantity represent?

$$\sum_j \mathbf{a}_j^2$$

This is the square of the magnitude of the vector!

The magnitude of a vector is called its **L_2 norm** and is denoted $\|\mathbf{a}\|_2$.

$$\|\mathbf{a}\|_2 = \sqrt{\sum_j \mathbf{a}_j^2}$$

Therefore, our above sum is $\|\mathbf{a}\|_2^2$.

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2$$

Think of $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots]$ as a vector.

What does the following quantity represent?

$$\sum_j \mathbf{a}_j^2$$

This is the square of the magnitude of the vector!

The magnitude of a vector is called its **L_2 norm** and is denoted $\|\mathbf{a}\|_2$.

$$\|\mathbf{a}\|_2 = \sqrt{\sum_j \mathbf{a}_j^2}$$

Therefore, our above sum is $\|\mathbf{a}\|_2^2$.

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2 \leq \frac{\|\mathbf{a}\|_2^2}{w}$$

Think of $[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \dots]$ as a vector.

What does the following quantity represent?

$$\sum_j \mathbf{a}_j^2$$

Great exercise: Prove that the L_2 norm of a vector is never greater than the L_1 norm.

This is the square of the magnitude of the vector.
The magnitude of a vector is often denoted $\|\mathbf{a}\|_2$.

$$\|\mathbf{a}\|_2 = \sqrt{\sum_j \mathbf{a}_j^2}$$

Therefore, our above sum is $\|\mathbf{a}\|_2^2$.

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{1}{w} \sum_{j \neq i} \mathbf{a}_j^2 \leq \frac{\|\mathbf{a}\|_2^2}{w}$$

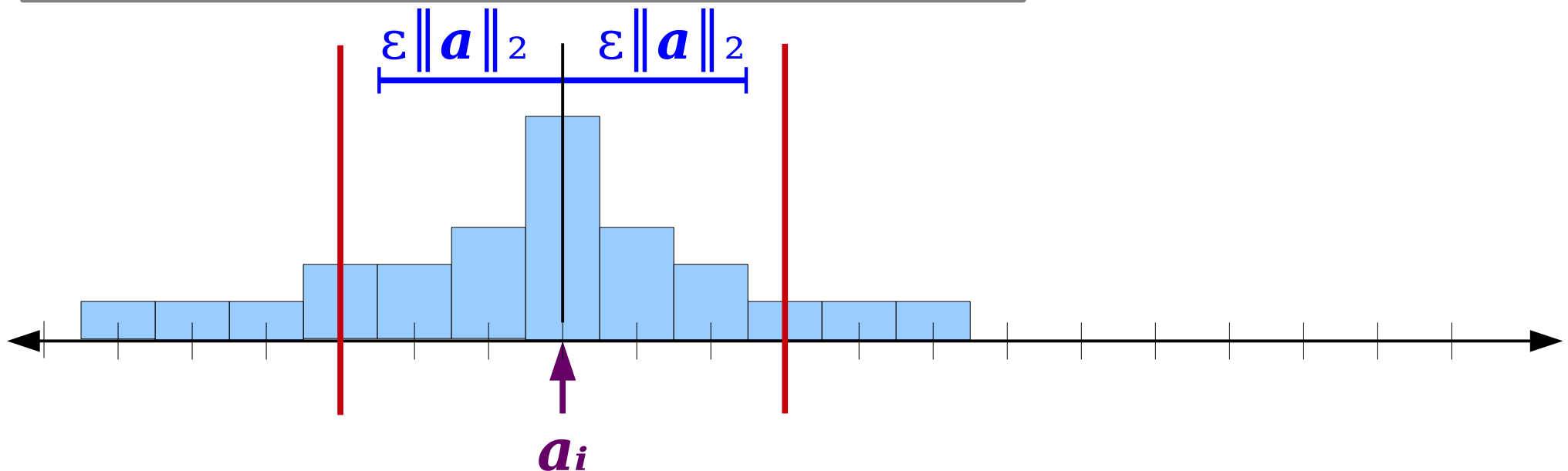
Goal: Make an estimator $\hat{\mathbf{a}}$ for some quantity \mathbf{a} where

With probability at least $1 - \delta$,

$$|\hat{\mathbf{a}} - \mathbf{a}| \leq \varepsilon \cdot \text{size}(\text{input})$$

Probably
Approximately Correct

for some measure of the size of the input.



$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{\|\mathbf{a}\|_2^2}{w}$$

$$\Pr[|\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i| > \varepsilon \|\boldsymbol{a}\|_2]$$

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2]$$

Chebyshev's inequality says that

$$\Pr[\|X - \mathbb{E}[X]\| \geq c] \leq \frac{\text{Var}[X]}{c^2}.$$

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2] \\ \leq \frac{\text{Var}[\hat{\mathbf{a}}_i]}{(\varepsilon \|\mathbf{a}\|_2)^2}$$

Chebyshev's inequality says that

$$\Pr[\|X - \mathbb{E}[X]\| \geq c] \leq \frac{\text{Var}[X]}{c^2}.$$

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2] \\ \leq \frac{\text{Var}[\hat{\mathbf{a}}_i]}{(\varepsilon \|\mathbf{a}\|_2)^2}$$

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{\|\mathbf{a}\|_2^2}{w}$$

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2]$$

$$\leq \frac{\text{Var}[\hat{\mathbf{a}}_i]}{(\varepsilon \|\mathbf{a}\|_2)^2}$$

$$\leq \frac{\|\mathbf{a}\|_2^2}{w} \cdot \frac{1}{(\varepsilon \|\mathbf{a}\|_2)^2}$$

$$\text{Var}[\hat{\mathbf{a}}_i] \leq \frac{\|\mathbf{a}\|_2^2}{w}$$

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2]$$

$$\leq \frac{\text{Var}[\hat{\mathbf{a}}_i]}{(\varepsilon \|\mathbf{a}\|_2)^2}$$

$$\leq \frac{\|\mathbf{a}\|_2^2}{w} \cdot \frac{1}{(\varepsilon \|\mathbf{a}\|_2)^2}$$

$$= \frac{1}{w \varepsilon^2}$$

Goal: Make an estimator $\hat{\mathbf{a}}$ for some quantity \mathbf{a} where

With probability at least $1 - \delta$,
 $|\hat{\mathbf{a}} - \mathbf{a}| \leq \varepsilon \cdot \text{size}(\text{input})$

Probably
Approximately Correct

for some measure of input size.

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2] \leq \frac{1}{w \varepsilon^2}$$

Pick $w = 4 \cdot \varepsilon^{-2}$. Then

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2] \leq \frac{1}{4}.$$

We now have a single estimator with a not-so-great chance of giving a good estimate.

How do we fix this?

How to Build an Estimator

	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	2-independence breaks up products; ± 1 variables have zero expected value.
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	Two-sided error; compute variance and use Chebyshev's inequality.
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

How to Build an Estimator

	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	2-independence breaks up products; ± 1 variables have zero expected value.
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	Two-sided error; compute variance and use Chebyshev's inequality.
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	

Running in Parallel

- Let's suppose that we run ***d*** independent copies of this data structure. Each has its own independently randomly chosen hash function.
- To ***increment***(x) in the overall structure, we call ***increment***(x) on each of the underlying data structures.
- The probability that at least one of them provides a good estimate is high.
- ***Question:*** How do you know which one?

Estimator 1:
137

Estimator 2:
271

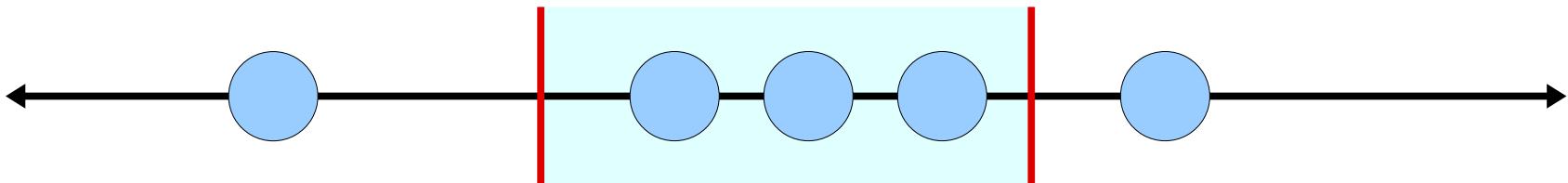
Estimator 3:
166

Estimator 4:
103

Estimator 5:
261

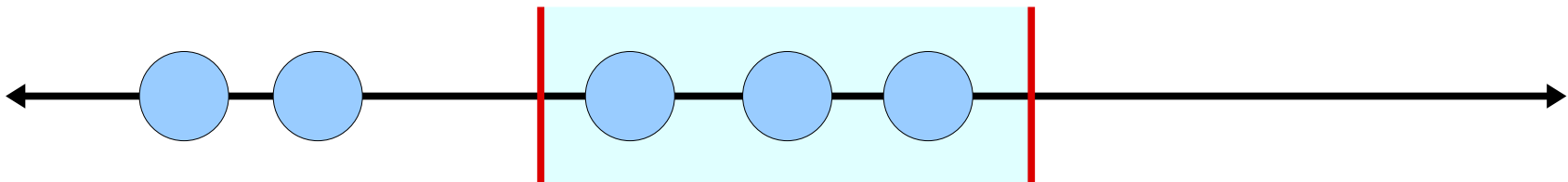
Working with the Median

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving the right answer.
- **Intuition:** The only way that the median isn't in the “good” area is if **at least half** the estimates are in the “bad” area.
- Each individual data structure has a “reasonable” chance to be good, so this is very unlikely.



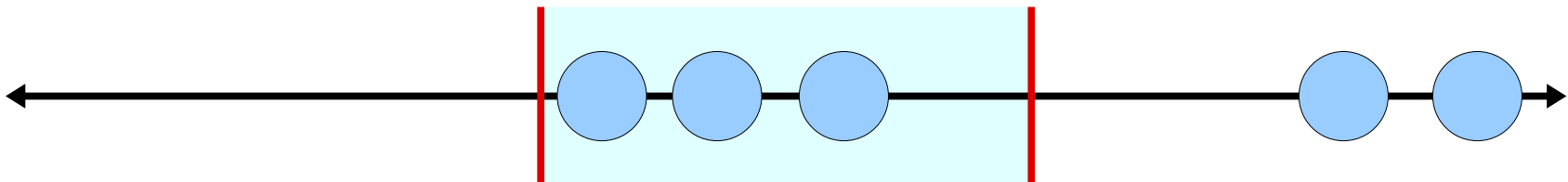
Working with the Median

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving the right answer.
- **Intuition:** The only way that the median isn't in the “good” area is if **at least half** the estimates are in the “bad” area.
- Each individual data structure has a “reasonable” chance to be good, so this is very unlikely.



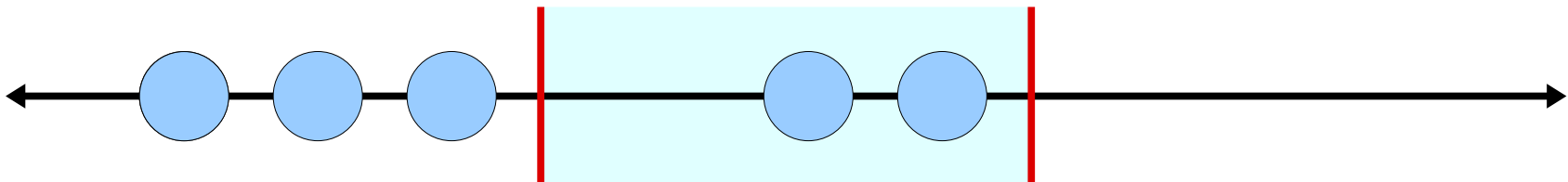
Working with the Median

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving the right answer.
- **Intuition:** The only way that the median isn't in the “good” area is if **at least half** the estimates are in the “bad” area.
- Each individual data structure has a “reasonable” chance to be good, so this is very unlikely.



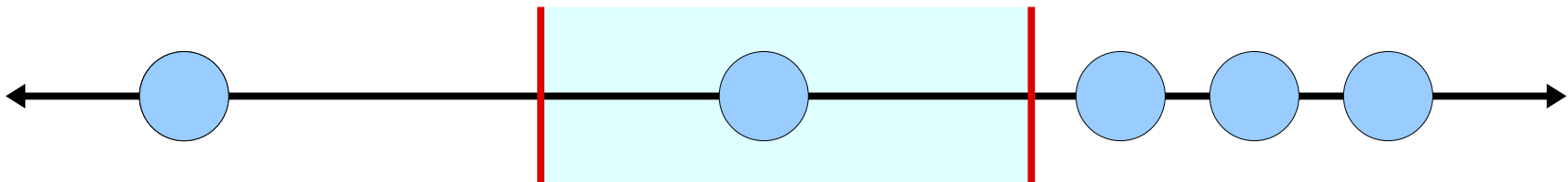
Working with the Median

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving the right answer.
- **Intuition:** The only way that the median isn't in the “good” area is if **at least half** the estimates are in the “bad” area.
- Each individual data structure has a “reasonable” chance to be good, so this is very unlikely.



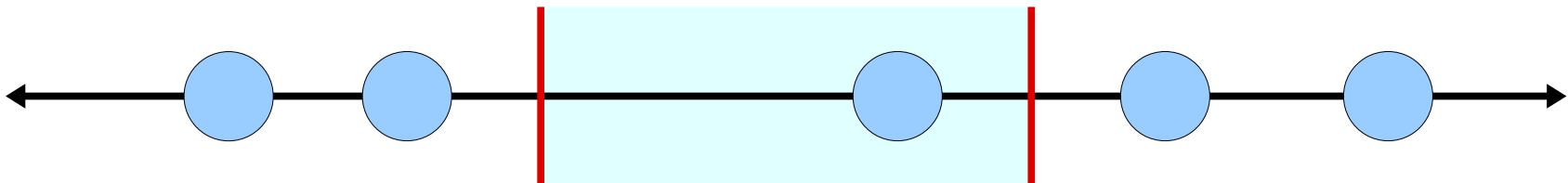
Working with the Median

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving the right answer.
- **Intuition:** The only way that the median isn't in the “good” area is if **at least half** the estimates are in the “bad” area.
- Each individual data structure has a “reasonable” chance to be good, so this is very unlikely.



Working with the Median

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving the right answer.
- **Intuition:** The only way that the median isn't in the “good” area is if **at least half** the estimates are in the “bad” area.
- Each individual data structure has a “reasonable” chance to be good, so this is very unlikely.



The Setup

- Let D denote a random variable equal to the number of data structures that produce an answer *not* within $\varepsilon ||\mathbf{a}||_2$ of the true answer.
- Since each independent data structure has failure probability at most $1/4$, we can upper-bound D with a $\text{Binom}(d, 1/4)$ variable.
- We want to know $\Pr[D > d / 2]$.
- How can we determine this?

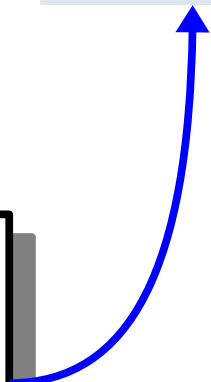
Chernoff Bounds

- The **Chernoff bound** says that if $X \sim \text{Binom}(n, p)$ and $p < 1/2$, then

$$\Pr[X \geq n/2] < e^{-n \cdot z(p)}$$

where $z(p) = (1/2 - p)^2 / 2p$.

Intuition: For any fixed value of p , this quantity decays exponentially quickly as a function of n . It's extremely unlikely that more than half our estimates will be bad.



Chernoff Bounds

- The **Chernoff bound** says that if $X \sim \text{Binom}(n, p)$ and $p < 1/2$, then

$$\Pr[X \geq n/2] < e^{-n \cdot z(p)}$$

where $z(p) = (1/2 - p)^2 / 2p$.

- In our case, $D \sim \text{Binom}(d, 1/4)$, so we know that

$$\Pr[D \geq \frac{d}{2}] \leq e^{-n \cdot z(1/4)} = e^{-d/8}$$

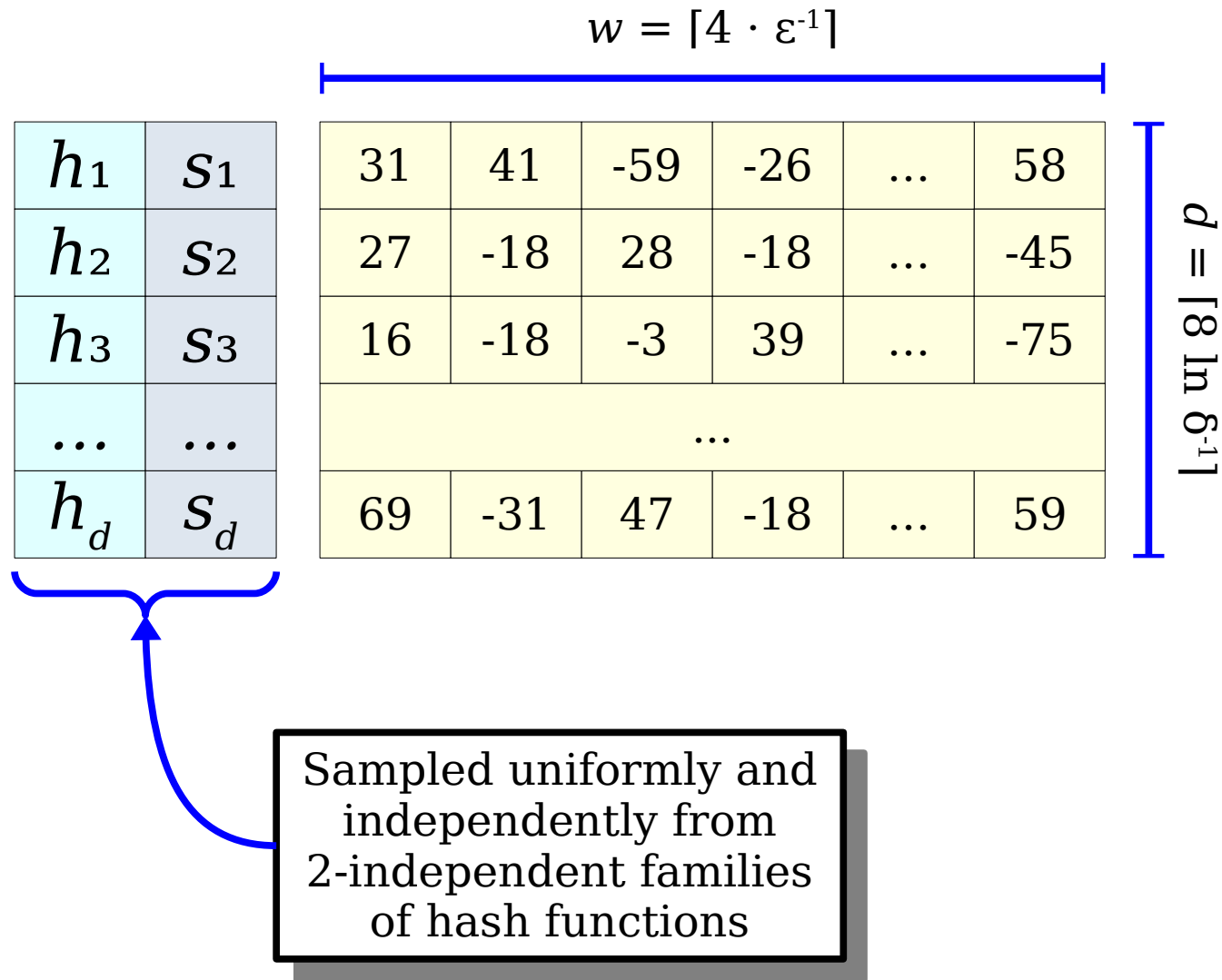
- Therefore, choosing **$d = 8 \log \delta^{-1}$** ensures that

$$\Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2] \leq \Pr[D \geq \frac{d}{2}] \leq \delta$$

How to Build an Estimator



	<i>Count-Min Sketch</i>	<i>Count Sketch</i>
Step One: Build a Simple Estimator	Hash items to counters; add +1 when item seen.	Hash items to counters; add ± 1 when item seen.
Step Two: Compute Expected Value of Estimator	Sum of indicators; 2-independent hashes have low collision rate.	2-independence breaks up products; ± 1 variables have zero expected value.
Step Three: Apply Concentration Inequality	One-sided error; use expected value and Markov's inequality.	Two-sided error; compute variance and use Chebyshev's inequality.
Step Four: Replicate to Boost Confidence	Take min; only fails if all estimates are bad.	Take median; only can fail if half of estimates are wrong; use Chernoff.

The Count Sketch



The Count Sketch

$$w = \lceil 4 \cdot \varepsilon^{-1} \rceil$$

							
h_1	s_1	31	41	-59	-26	...	58
h_2	s_2	27	-18	28	-18	...	-45
h_3	s_3	16	-18	-3	39	...	-75
...					
h_d	s_d	69	-31	47	-18	...	59
							
		$d = \lceil 8 \ln 8 \cdot \varepsilon^{-1} \rceil$					

```

increment(x):
  for i = 1 ... d:
    count[i][hi(x)] += si(x)
    
```

The Count Sketch

$$w = \lceil 4 \cdot \varepsilon^{-1} \rceil$$

							
h_1	s_1	31	41	-59	-26	...	58
h_2	s_2	27	-18	28	-18	...	-45
h_3	s_3	16	-18	-3	39	...	-75
...					
h_d	s_d	69	-31	47	-18	...	59
							
		$d = \lceil 8 \ln 8^{-1} \rceil$					

```

increment(x):
  for i = 1 ... d:
    count[i][hi(x)] += si(x)
    
```

The Count Sketch

$$w = \lceil 4 \cdot \varepsilon^{-1} \rceil$$



h_1	s_1	31	40	-59	-26	...	58
h_2	s_2	27	-18	28	-19	...	-45
h_3	s_3	16	-18	-3	40	...	-75
...					
h_d	s_d	69	-31	47	-18	...	58

$d = \lceil 8 \ln 6^{-1} \rceil$

```
increment(x):
    for i = 1 ... d:
        count[i][hi(x)] += si(x)
```

The Count Sketch

$$w = \lceil 4 \cdot \varepsilon^{-1} \rceil$$


							
h_1	s_1	31	40	-59	-26	...	58
h_2	s_2	27	-18	28	-19	...	-45
h_3	s_3	16	-18	-3	40	...	-75
...					
h_d	s_d	69	-31	47	-18	...	58
							
		$d = \lceil 8 \ln 8 \cdot \varepsilon^{-1} \rceil$					

```

increment(x):
  for i = 1 ... d:
    count[i][hi(x)] += si(x)
    
```

The Count Sketch

$$w = \lceil 4 \cdot \varepsilon^{-1} \rceil$$


							
h_1	s_1	31	40	-59	-26	...	58
h_2	s_2	27	-18	28	-19	...	-45
h_3	s_3	16	-18	-3	40	...	-75
...					
h_d	s_d	69	-31	47	-18	...	58
							
		$d = \lceil 8 \ln 8^{-1} \rceil$					

```
increment(x):
  for i = 1 ... d:
    count[i][hi(x)] += si(x)
```

```
estimate(x):
  options = []
  for i = 1 ... d:
    options += count[i][hi(x)] * si(x)
  return medianOf(options)
```

The Count Sketch

$$w = \lceil 4 \cdot \varepsilon^{-1} \rceil$$

							
h_1	s_1	31	40	-59	-26	...	58
h_2	s_2	27	-18	28	-19	...	-45
h_3	s_3	16	-18	-3	40	...	-75
...					
h_d	s_d	69	-31	47	-18	...	58
							
		$d = \lceil 8 \ln 8^{-1} \rceil$					

```
increment(x):
  for i = 1 ... d:
    count[i][hi(x)] += si(x)
```

```
estimate(x):
  options = []
  for i = 1 ... d:
    options += count[i][hi(x)] * si(x)
  return medianOf(options)
```


The Final Analysis

- Here's a comparison of these two structures.
- **Question to ponder:** When is a count-min sketch better than a count sketch, and vice-versa?

Count-Min Sketch

Space: $\Theta(\varepsilon^{-1} \log \delta^{-1})$

increment: $\Theta(\log \delta^{-1})$

estimate: $\Theta(\log \delta^{-1})$

Accuracy: within $\varepsilon \|a\|_1$.

Count Sketch

Space: $\Theta(\varepsilon^{-2} \log \delta^{-1})$

increment: $\Theta(\log \delta^{-1})$

estimate: $\Theta(\log \delta^{-1})$

Accuracy: within $\varepsilon \|a\|_2$

Bloom Filters

Exact Membership Queries

- Suppose you're in a memory-constrained environment where every bit of memory counts, and you want to store a set of n items drawn from some universe U .
- Examples:
 - You're working on an embedded device with some maximum amount of working RAM.
 - You're working with large n (say, $n = 10^9$) on a modern machine.
 - You're building a consumer application like a web browser and don't want to hog all system resources.
- **Question:** How many bits of memory are needed to do this?

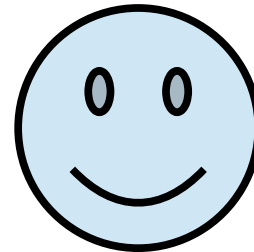
A Quick Detour

Goal: Design a simple data structure that can hold a single one of the objects shown to the right.

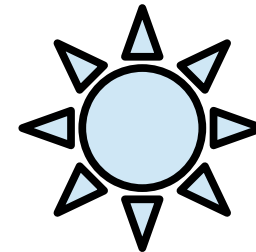
What is the minimum number of *bits* (not *words*) required for this data structure in the worst case?

We can get away with four bits by numbering each item and just storing the number.

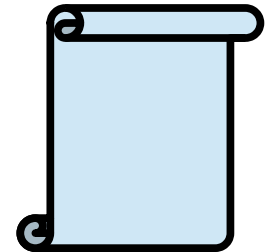
Question: Can we do better?



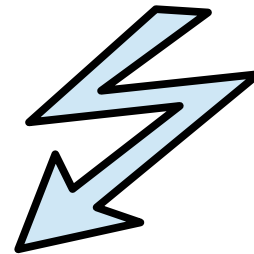
0000



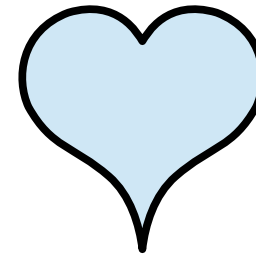
0001



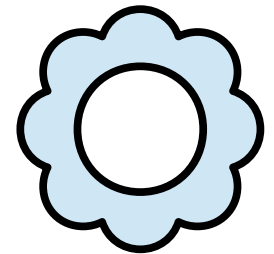
0010



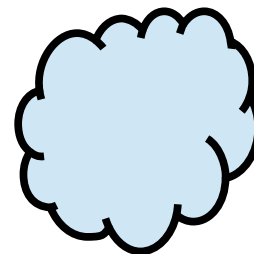
0011



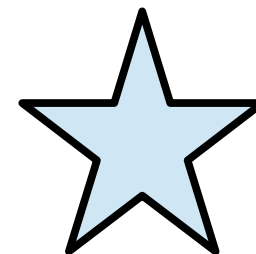
0100



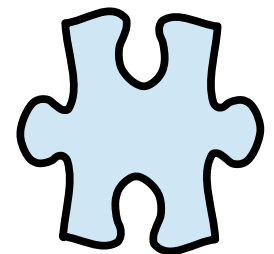
0101



0110



0111

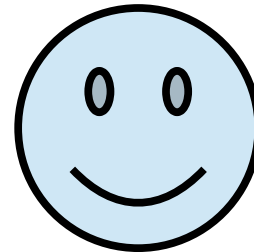


1000

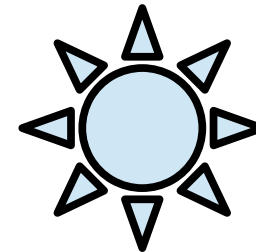
Goal: Design a simple data structure that can hold a single one of the objects shown to the right.

Claim: Every data structure for this problem must use at least four bits of memory in the worst case.

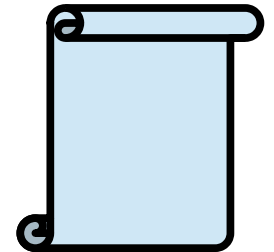
Proof: If we always use three or fewer bits, there are at most $2^3 = 8$ combinations of those bits, not enough to uniquely identify one of the nine different items.



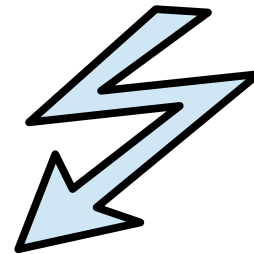
00000



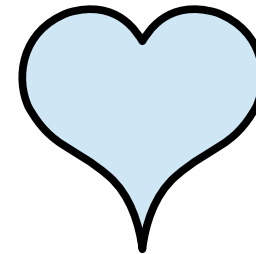
00001



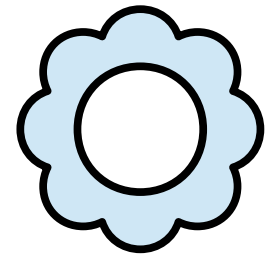
00010



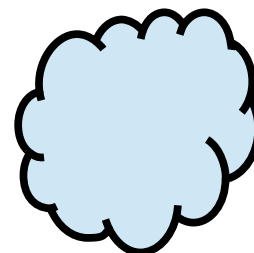
00011



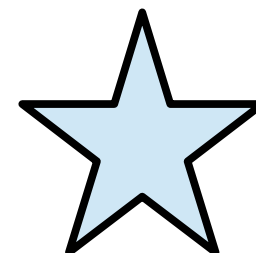
00100



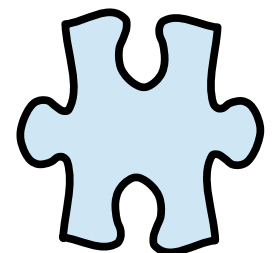
00101



00110



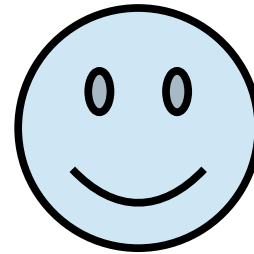
00111



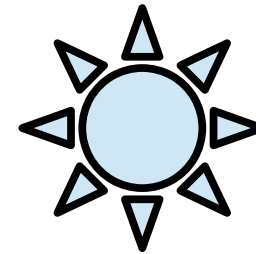
01000

Theorem: A data structure that stores one object out of a set of k possibilities must use at least $\lg k$ bits in the worst case.

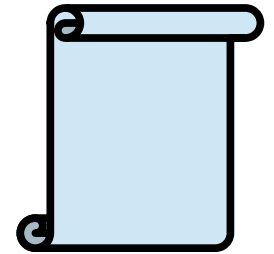
Proof: Using fewer than $\lg k$ bits means there are fewer than $2^{\lg k} = k$ possible combinations of those bits, not enough to uniquely identify each item out of the set.



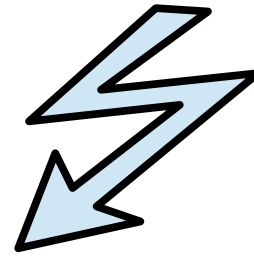
00000



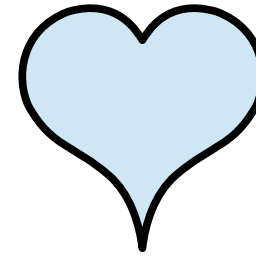
00001



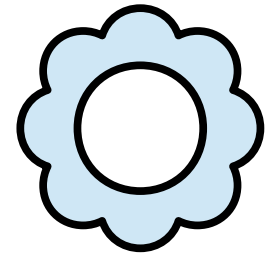
00100



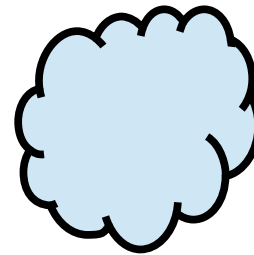
00111



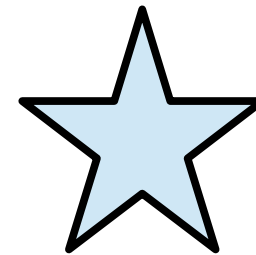
01100



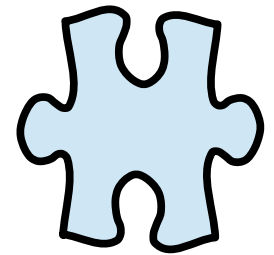
01011



01110



01111



10000

Question: How much memory is needed to solve the exact membership query problem?

Suppose we want to store a set $S \subseteq U$ of size n . How many bits of memory do we need?

Number of n -element subsets of universe U :

$$\binom{|U|}{n}$$

Bits needed:

$$\Omega(n \lg |U| - n \lg n)$$

$$\begin{aligned} & \lg \binom{|U|}{n} \\ &= \lg \left(\frac{|U|!}{n! (|U| - n)!} \right) \\ &\geq \lg \left(\frac{(|U| - n)^n}{n^n} \right) \\ &= n \lg \left(\frac{|U| - n}{n} \right) \\ &= n \lg \left(\frac{|U|}{n} - 1 \right) \\ &\geq n \lg \left(\frac{|U|}{n} \right) - n \\ &\geq n \lg |U| - n \lg n - n \\ &= \Omega(n \lg |U| - n \lg n) \end{aligned}$$

Bitten by Bits

- Assuming $|U| \gg n$, we need $\Omega(n \log |U|)$ bits to encode a solution to the exact membership query problem.
- If we're resource-constrained, this might be way too many bits for us to fit things in memory.
 - Think $n = 10^8$ and U is the set of all possible URLs or human genomes.
- Can we do better?

Approximate Membership Queries

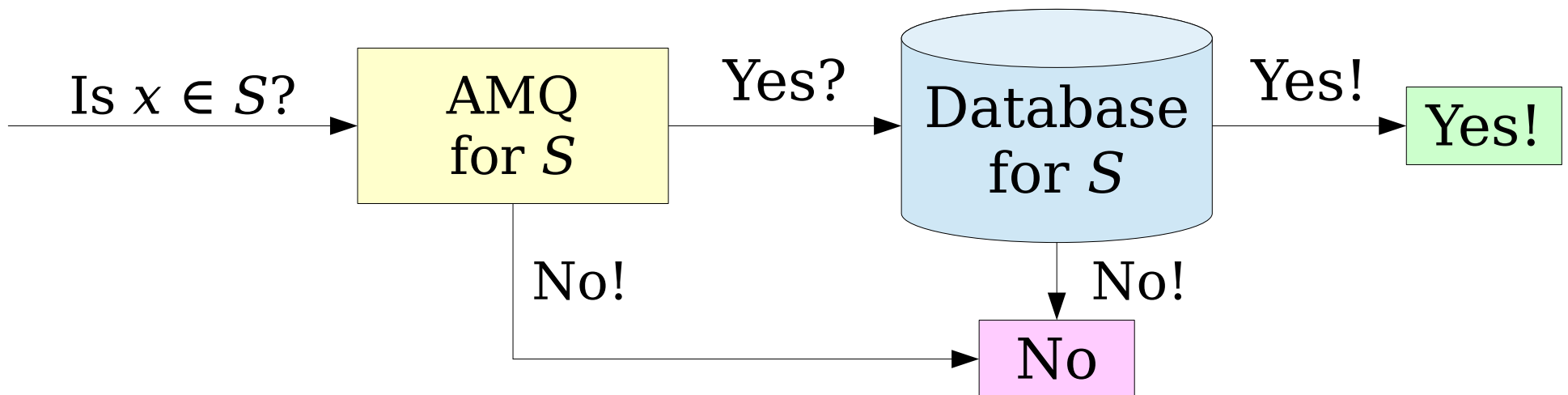
- The ***approximate membership query*** problem is the following:

Maintain a set S in a way that gives approximate answers to queries of the form “is $x \in S$?”

- Questions we need to answer:
 - How do you give an “approximate” answer to the question “is $x \in S$?”
 - Does this relaxation let us save memory?
- We’ll address each of these in turn.

Our Model

- **Goal:** Design our data structures to allow for false positives but not false negatives.
- That is:
 - if $x \in S$, we always return true, but
 - if $x \notin S$, we have a small probability of returning true.
- This is often a good idea in practice.



Our Model

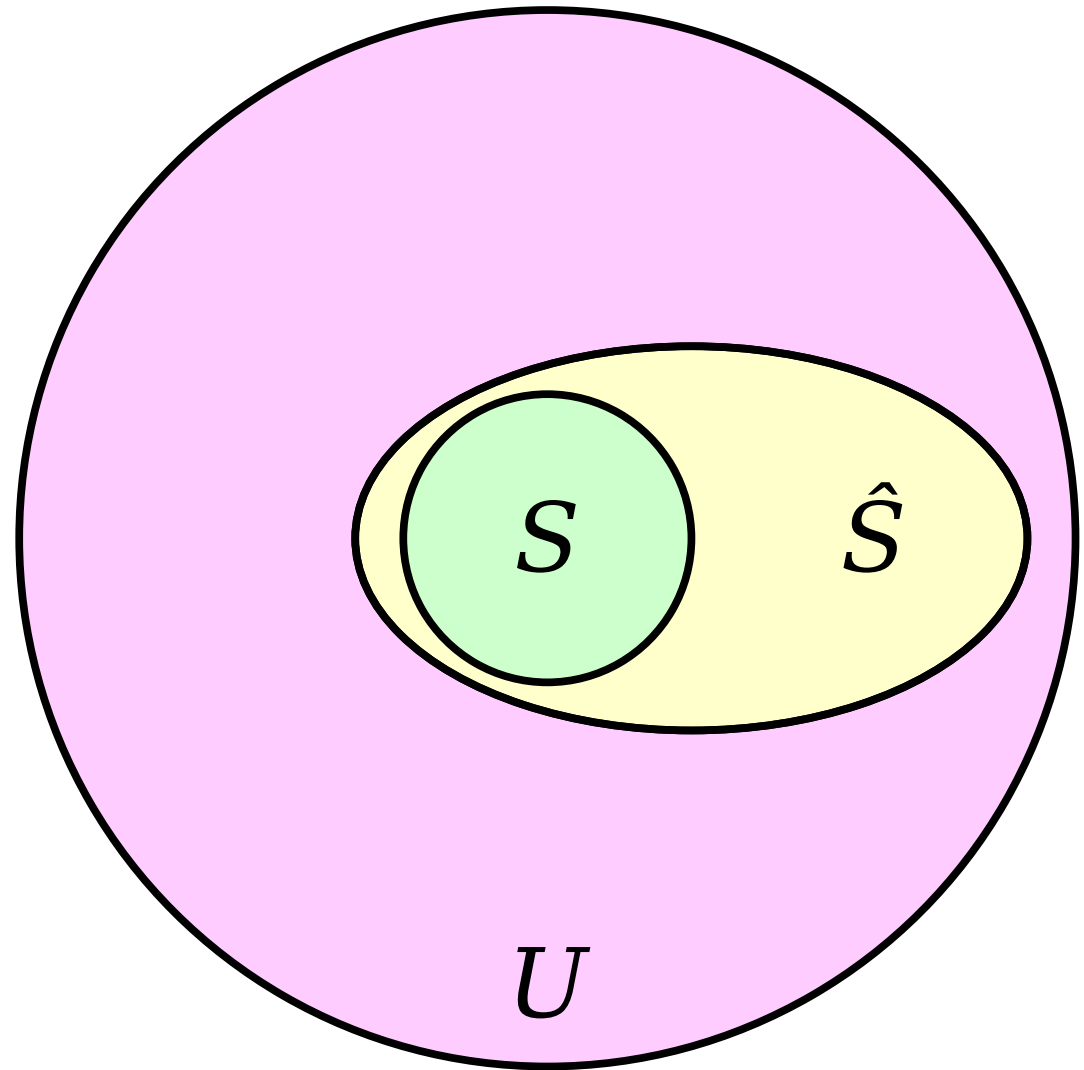
- Let's assume we have a tunable accuracy parameter $\varepsilon \in (0, 1)$.
- **Goal:** Design our data structure so that
 - if $x \in S$, we always return true;
 - if $x \notin S$, we return false with probability at least $1 - \varepsilon$; and
 - the amount of space we need depends only on n and ε , not on the size of the universe.
- Is this even possible?

Suppose we're storing
an n -element set S with
error rate ε .

How much memory is needed to solve
the approximate membership query problem?

Suppose we're storing an n -element set S with error rate ε .

Intuition: An AMQ structure stores a set \hat{S} : S plus approximately $\varepsilon|U|$ extra elements due to the error rate.

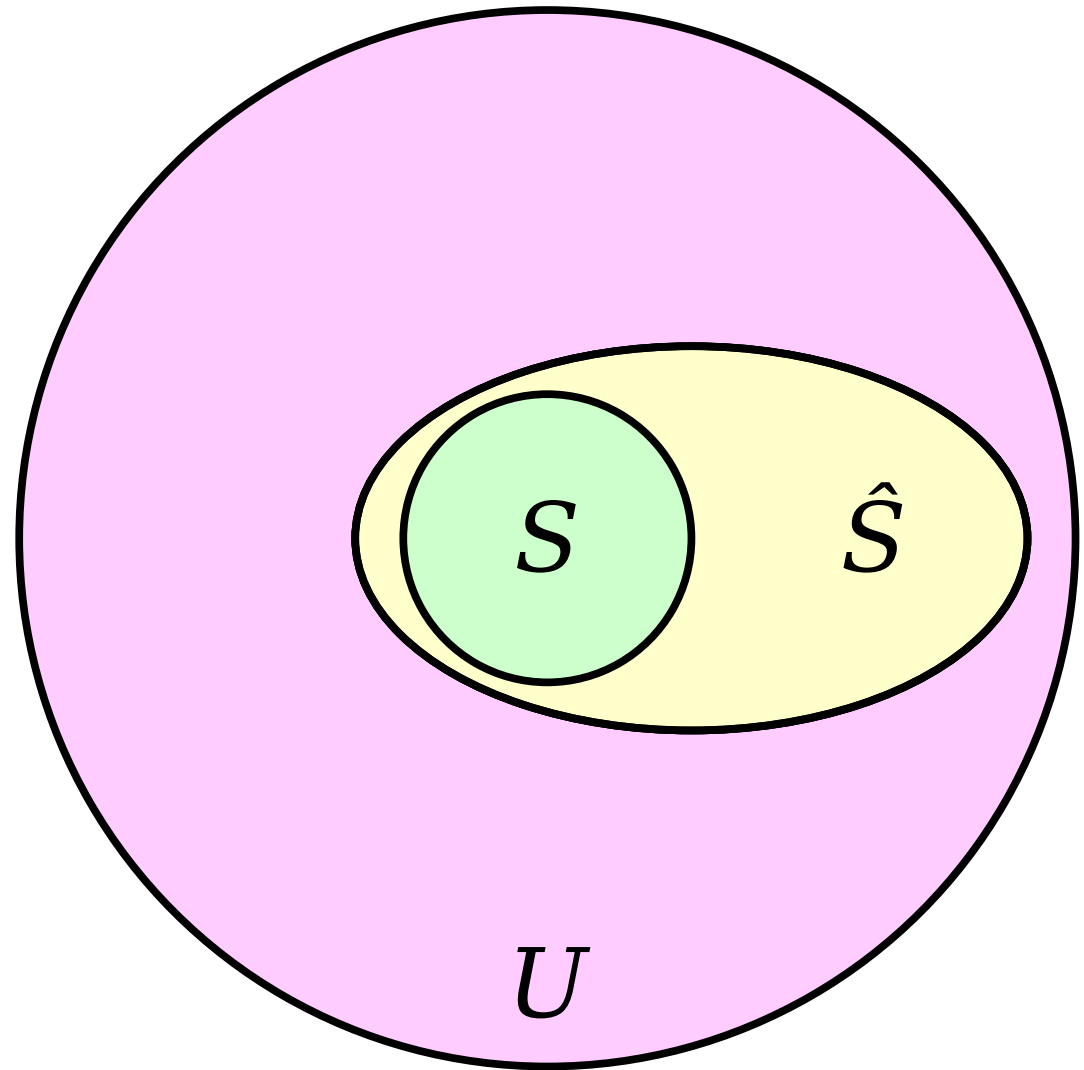


How much memory is needed to solve the approximate membership query problem?

Suppose we're storing an n -element set S with error rate ε .

Intuition: An AMQ structure stores a set \hat{S} : S plus approximately $\varepsilon|U|$ extra elements due to the error rate.

Importantly, we don't care *which* $\varepsilon|U|$ elements those are.

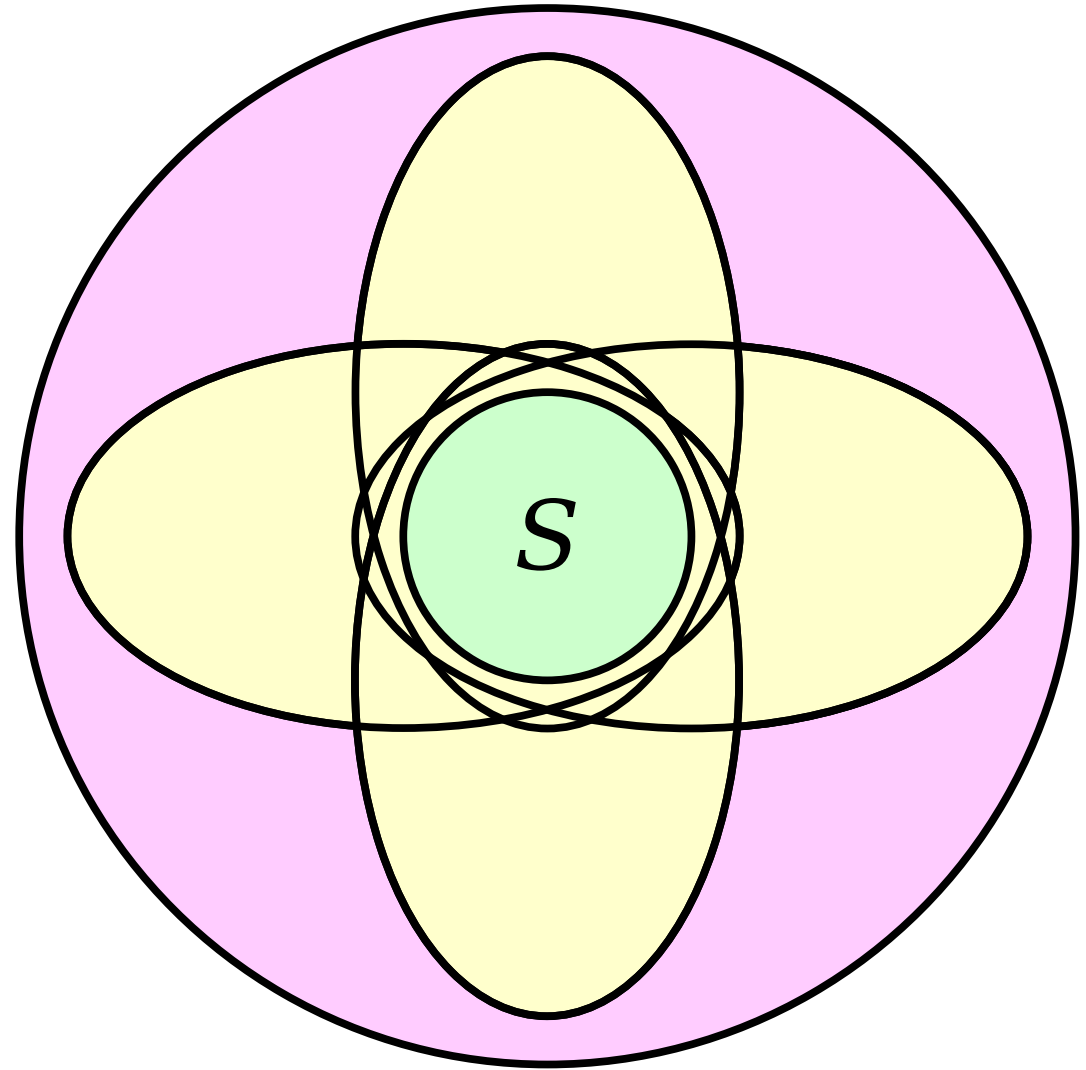


How much memory is needed to solve the approximate membership query problem?

Suppose we're storing an n -element set S with error rate ε .

Intuition: An AMQ structure stores a set \hat{S} : S plus approximately $\varepsilon|U|$ extra elements due to the error rate.

Importantly, we don't care *which* $\varepsilon|U|$ elements those are.

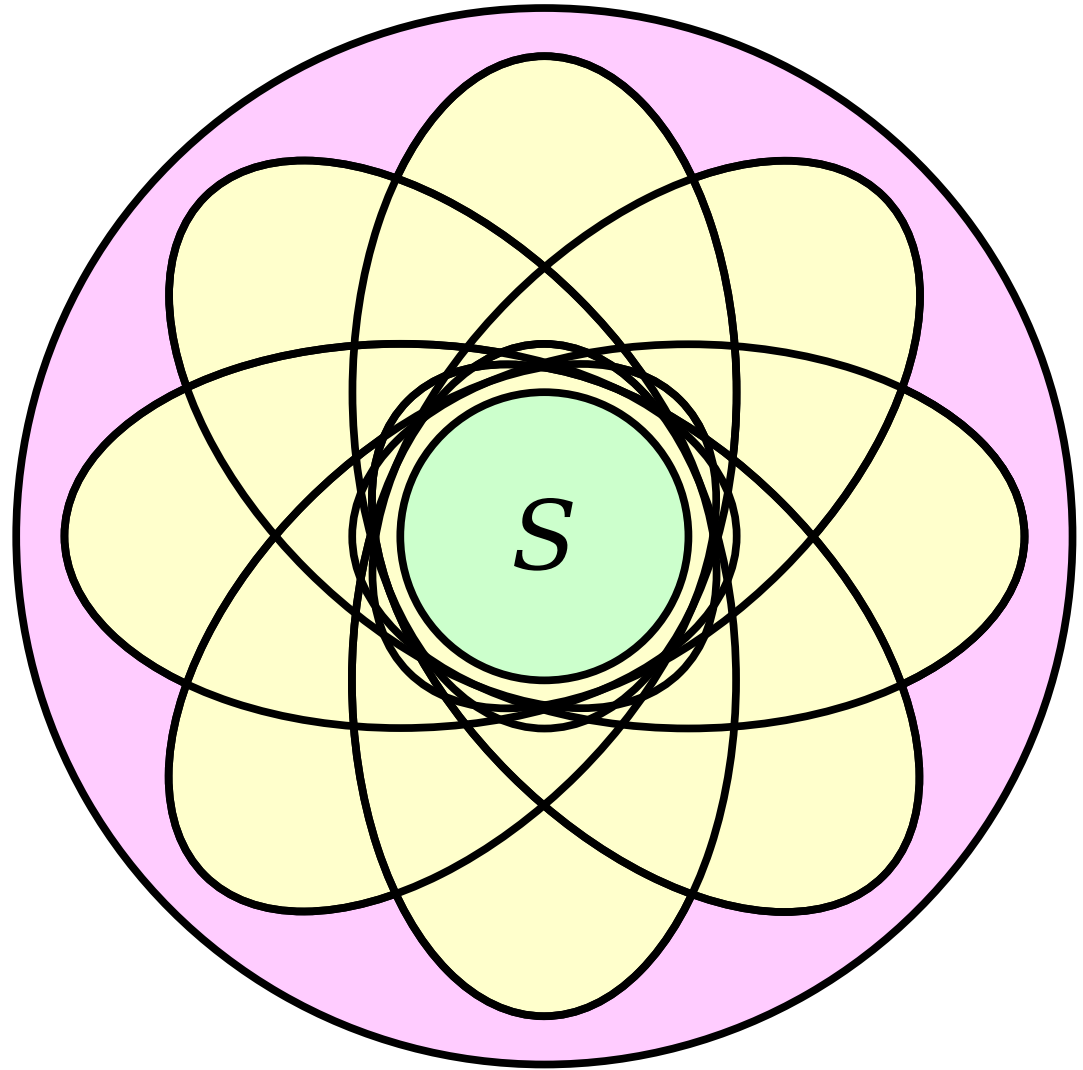


How much memory is needed to solve the approximate membership query problem?

Suppose we're storing an n -element set S with error rate ε .

Intuition: An AMQ structure stores a set \hat{S} : S plus approximately $\varepsilon|U|$ extra elements due to the error rate.

Importantly, we don't care *which* $\varepsilon|U|$ elements those are.



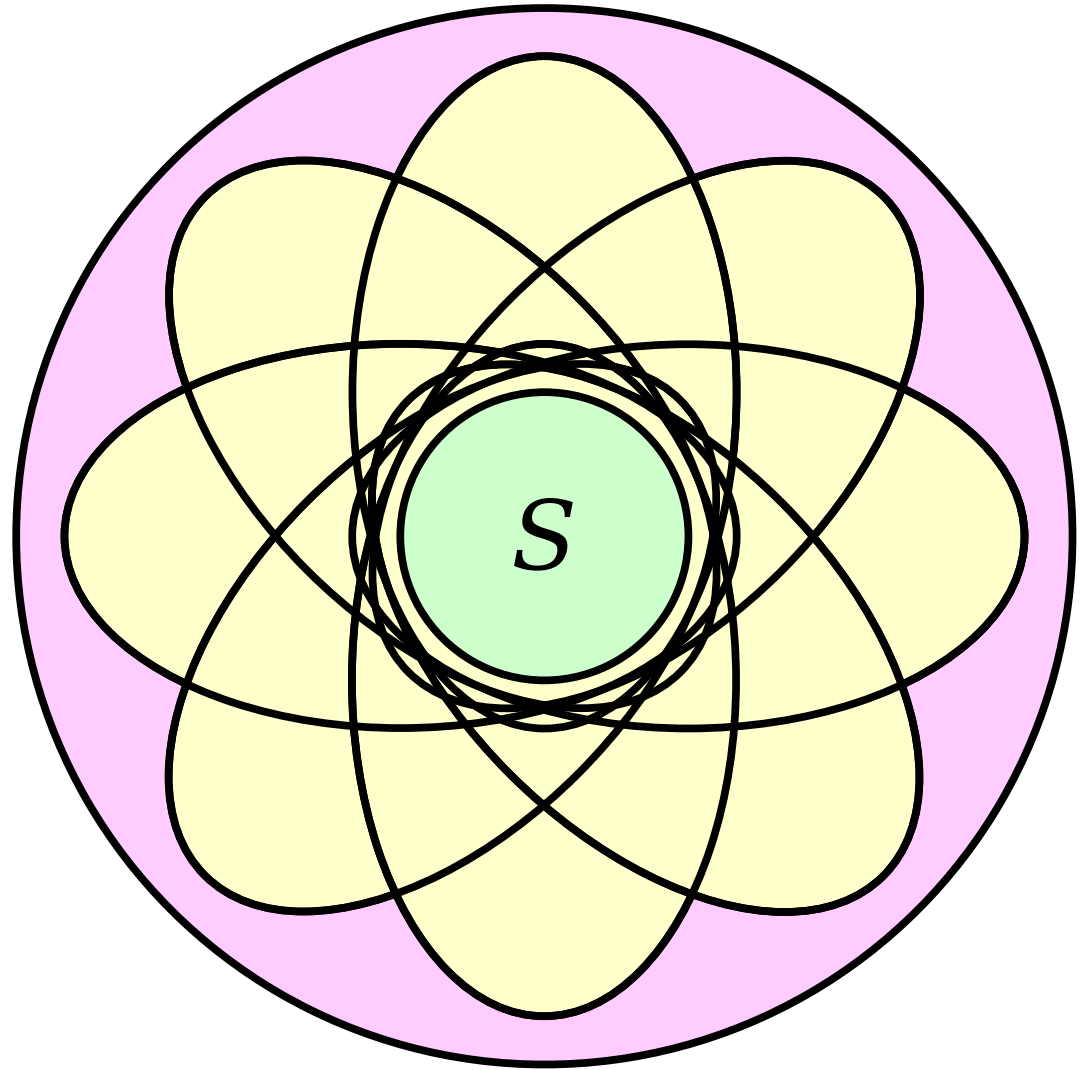
How much memory is needed to solve the approximate membership query problem?

Suppose we're storing an n -element set S with error rate ε .

Intuition: An AMQ structure stores a set \hat{S} : S plus approximately $\varepsilon|U|$ extra elements due to the error rate.

Importantly, we don't care *which* $\varepsilon|U|$ elements those are.

How does that affect our lower bound?



How much memory is needed to solve the approximate membership query problem?

How much memory is needed to solve
the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

How much memory is needed to solve the approximate membership query problem?

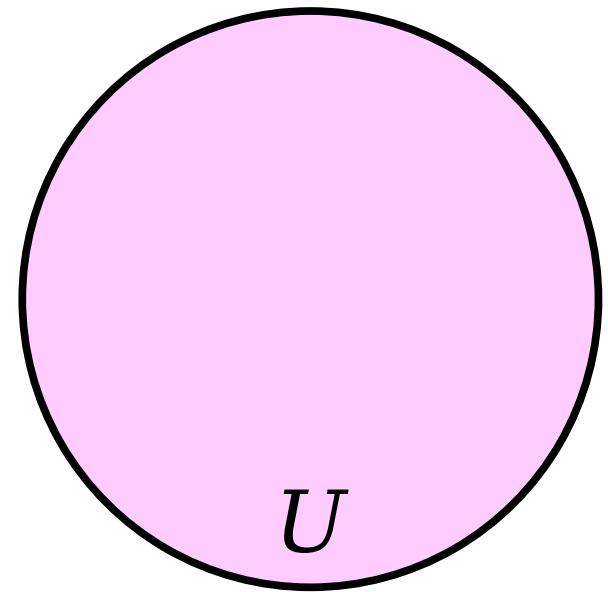
Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

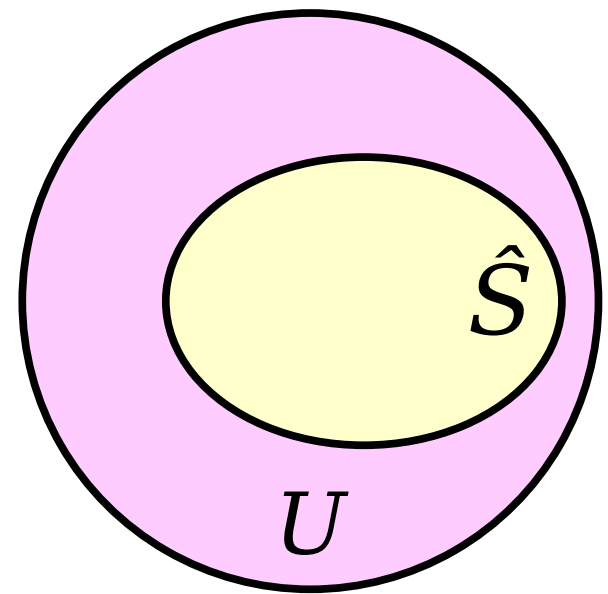
With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.



How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

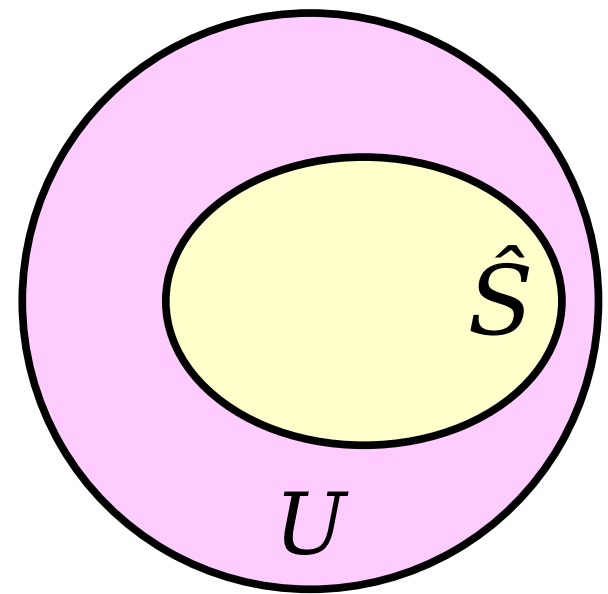


How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S .



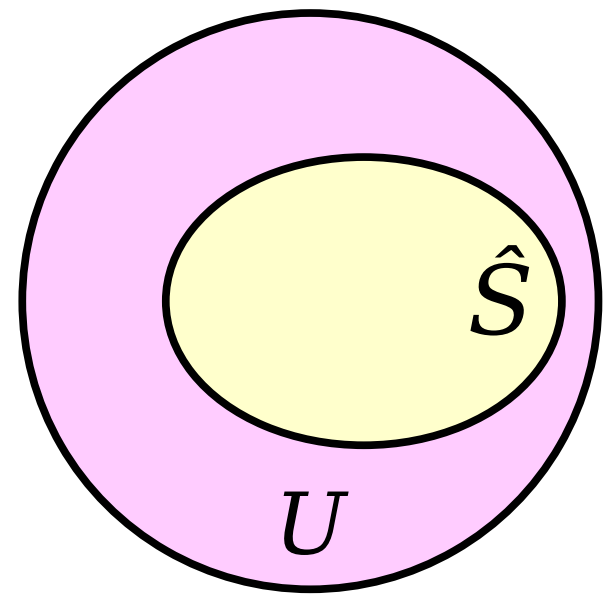
How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S . Bits needed:

$$\lg \binom{\varepsilon|U|}{n}$$



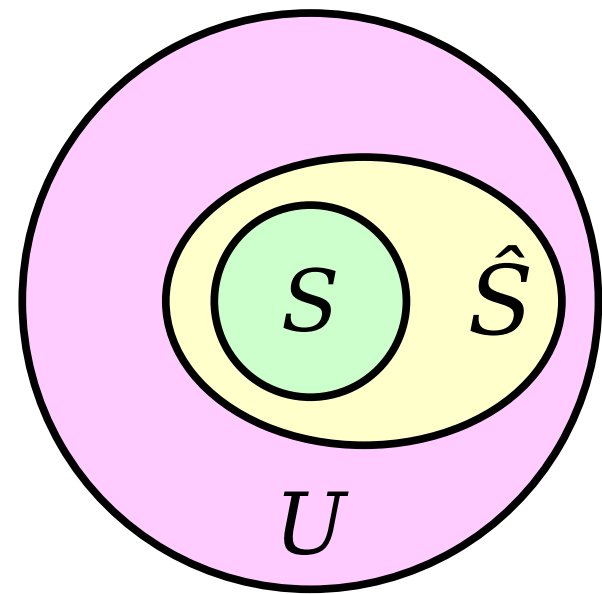
How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S . Bits needed:

$$\lg \binom{\varepsilon|U|}{n}$$



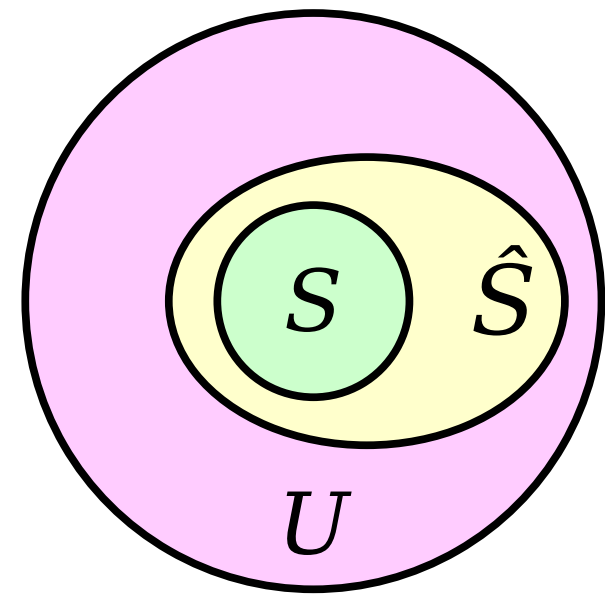
How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S . Bits needed:

$$\lg \binom{\varepsilon|U|}{n}$$



$$b + \lg \binom{\varepsilon|U|}{n}$$

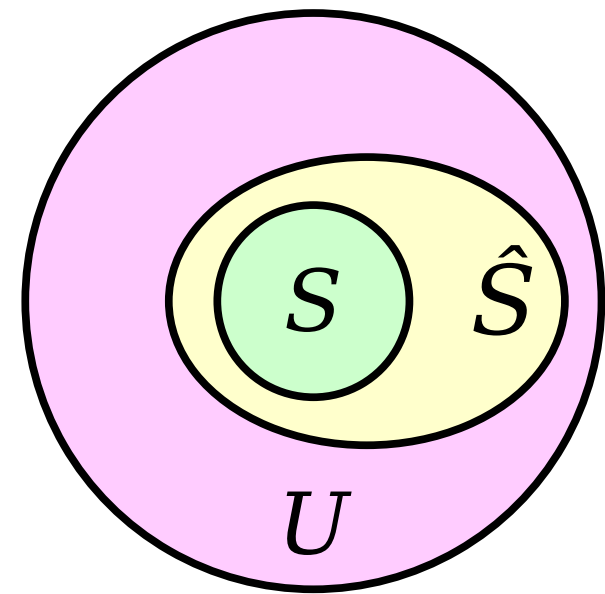
How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S . Bits needed:

$$\lg \binom{\varepsilon|U|}{n}$$



$$b + \lg \binom{\varepsilon|U|}{n}$$

(Number of bits needed to describe S .)

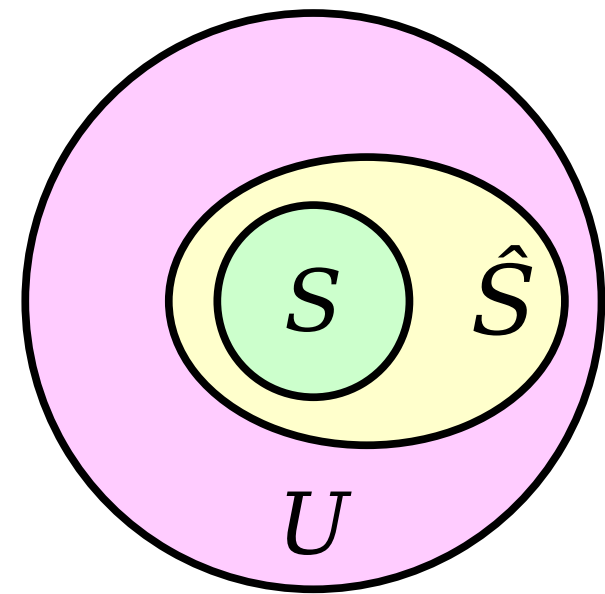
How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S . Bits needed:

$$\lg \binom{\varepsilon|U|}{n}$$



$$b + \lg \binom{\varepsilon|U|}{n} \geq \lg \binom{|U|}{n}$$

(Number of bits needed to describe S .)

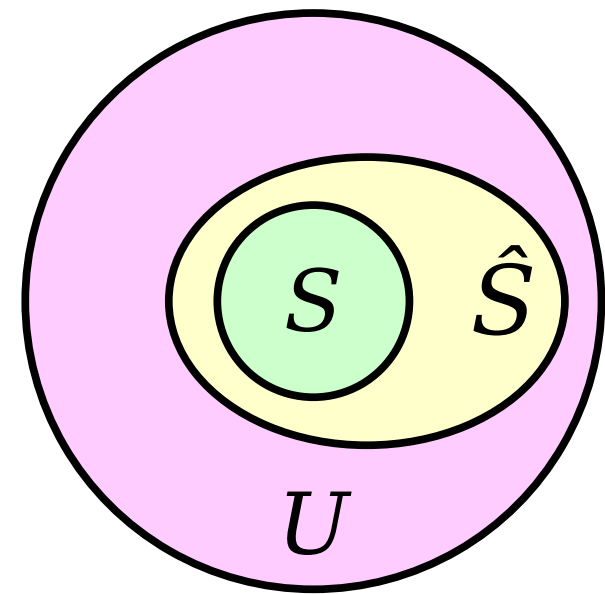
How much memory is needed to solve the approximate membership query problem?

Idea: We can use an AMQ, plus some more bits, to describe a set S .

With b bits, write down an AMQ structure. This describes a set $\hat{S} \subseteq U$ of size around $\varepsilon|U|$.

Write down some more bits to identify which n elements in \hat{S} make up the set S . Bits needed:

$$\lg \binom{\varepsilon|U|}{n}$$



$$b + \lg \binom{\varepsilon|U|}{n} \geq \lg \binom{|U|}{n}$$

(Number of bits needed to describe S) (Lower bound on number of bits needed.)

How much memory is needed to solve the approximate membership query problem?

Theorem: Assuming $\varepsilon|U| \gg n$, any AMQ structure needs at least approximately $n \lg \varepsilon^{-1}$ bits in the worst case.

This lower bound suggests we should aim for $\Theta(n \lg \varepsilon^{-1})$ bits of storage space.

How might we use them?

The math, if you're curious:

$$\begin{aligned} b &\geq \lg \left(\frac{\binom{|U|}{n}}{\binom{\varepsilon|U|}{n}} \right) \\ &\approx \lg \left(\frac{|U|^n / n!}{(\varepsilon|U|)^n / n!} \right) \\ &= \lg \varepsilon^{-n} \\ &= n \lg \varepsilon^{-1} \end{aligned}$$

How much memory is needed to solve the approximate membership query problem?

Bloom Filters

Idea 1: Adapt the “hash to a bucket” idea of the count-min and count sketches.

As an example, let's have

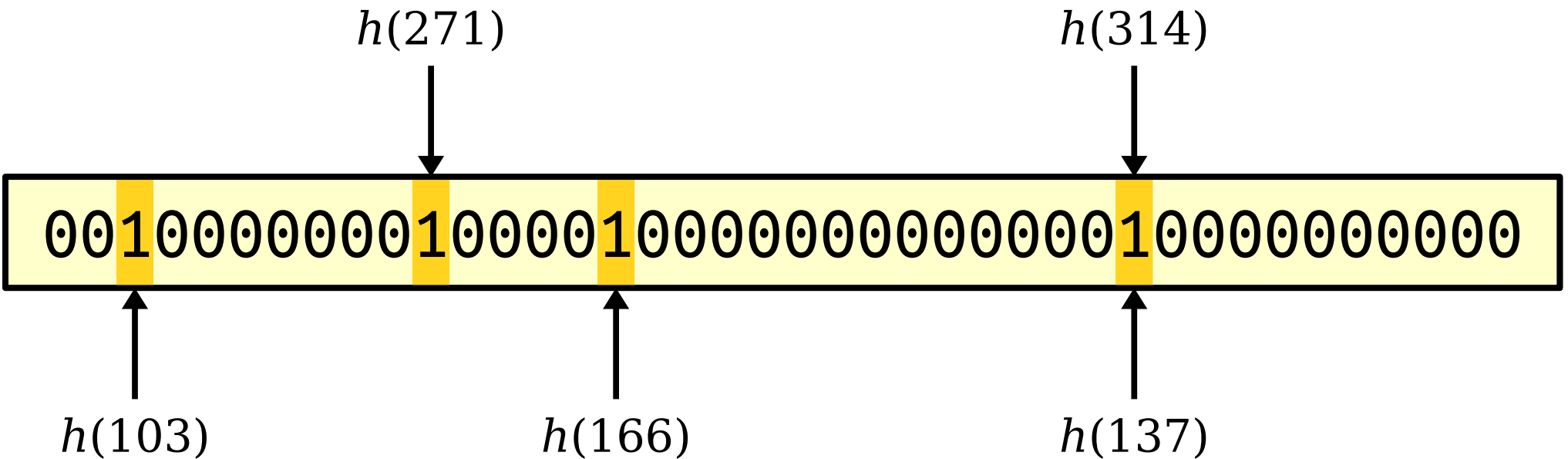
$$S = \{103, 137, 166, 271, 314\}$$
[illegible]

Number of bits: **$cn \lg \varepsilon^{-1}$**
(We'll pick c later.)

How can we approximate a set in a small number of bits and with a low error rate?

Idea 1: Adapt the “hash to a bucket” idea of the count-min and count sketches.

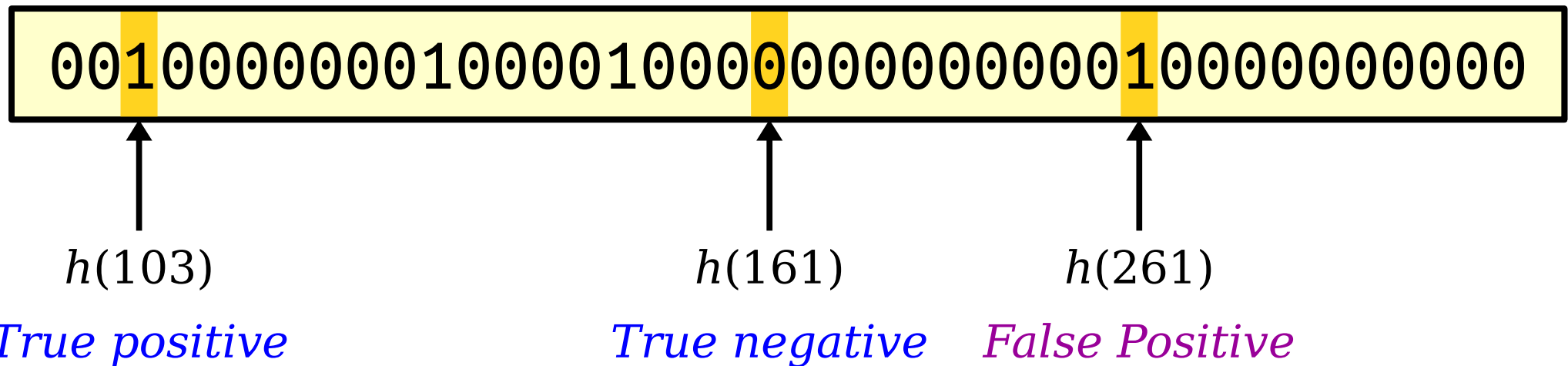
As an example, let's have

$$S = \{103, 137, 166, 271, 314\}$$


How can we approximate a set in a small number of bits and with a low error rate?

Idea 1: Adapt the “hash to a bucket” idea of the count-min and count sketches.

As an example, let's have
 $S = \{103, 137, 166, 271, 314\}$



How can we approximate a set in a small number of bits and with a low error rate?

Idea 1: Adapt the “hash to a bucket” idea of the count-min and count sketches.

Suppose we store a set of n elements in collection of $cn \lg \varepsilon^{-1}$ bits.

Question: What is our false positive rate?

00100000001000010000000000000000100000000000

Intuition: At most n of our $cn \lg \varepsilon^{-1}$ bits will be 1. We only have false positives if we see a 1. So our false positive rate will be about $1 / c \lg \varepsilon^{-1}$. That's not great.

How can we approximate a set in a small number of bits and with a low error rate?

Idea 2: Adapt the
“run in parallel”
approach of the
count-min sketch.

Make several copies of
the previous data
structure, each with a
random hash function.

000100000000000000000000000000001100000000010000000000001

001000000000100000100000000000000000000000000000000000000

0000001000000000000000001000000001000000100000000100

Can we get the same accuracy
while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

Question: Each copy provides its own estimate. Which one should we pick?

[illegible]

00100000001000010000 element to be in our set S . 0000

00000100000000000001000000100000100000000100

This entry is 0, so it's not possible for this element to be in our set S .

Can we get the same accuracy
while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

Question: Each copy provides its own estimate. Which one should we pick?

000100000000000000000000000000001100000000100000000000001

00100000000010000010000000000000000000000000100000000000000

0000001000000000000000100000000010000000100000000100

Can we get the same accuracy while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

Question: Each copy provides its own estimate. Which one should we pick?

0001000000000000000000000011000000001000000000001

001000000010000100000000000000001**00000000000000**

00000100000000000001000000010000010000000100

Can we get the same accuracy
while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

Question: Each copy provides its own estimate. Which one should we pick?

```
000100000000000000000000000011000000001000000000001
```

00100000001000010000000000000000100000000000

00000100000000000001000000010000

We only say
“yes” if all
bits are 1’s.

Can we get the same accuracy
while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

We have some fixed number of bits to use. How should we split them across these copies?

0010000000010000100000000000000000000000010000000000000

Can we get the same accuracy while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

We have some fixed number of bits to use. How should we split them across these copies?

0010010000100001000010

1001100001000000001000

Can we get the same accuracy while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

We have some fixed number of bits to use. How should we split them across these copies?

0010001100100100

0101000010001010

0000011100000010

Can we get the same accuracy while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

We have some fixed number of bits to use. How should we split them across these copies?

001000110010

001000110010

010100001000

000001110000

More copies means fewer bits per copy. The error rate per copy is higher, but there's more copies.

Can we get the same accuracy while using fewer bits?

Idea 2: Adapt the “run in parallel” approach of the count-min sketch.

Approach: Use one giant array. Have all hash functions edit and read that array.

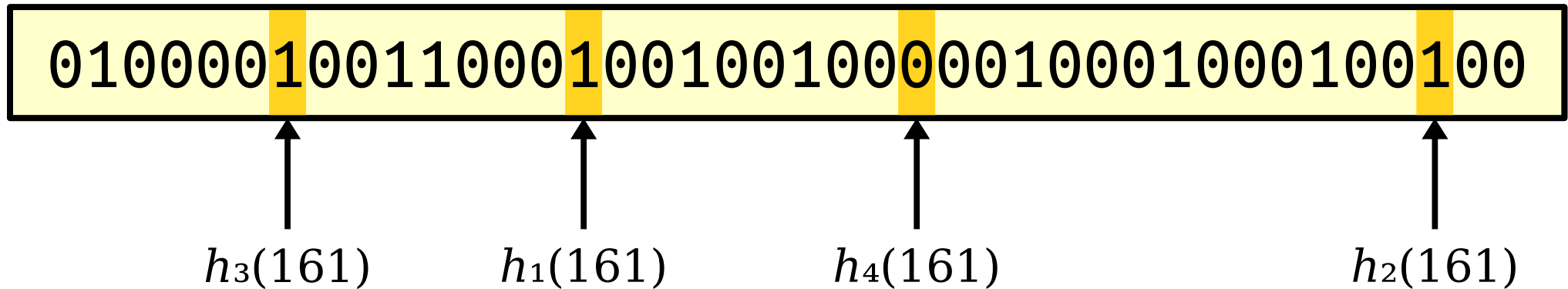
This is called a **Bloom filter**, named after its inventor.

01000010011000100100100000010001000100100

Can we get the same accuracy
while using fewer bits?

Assume we use k hash functions, each of which is chosen independently of the others. We'll pick k later on.

(In this example,
 $k = 4$.)

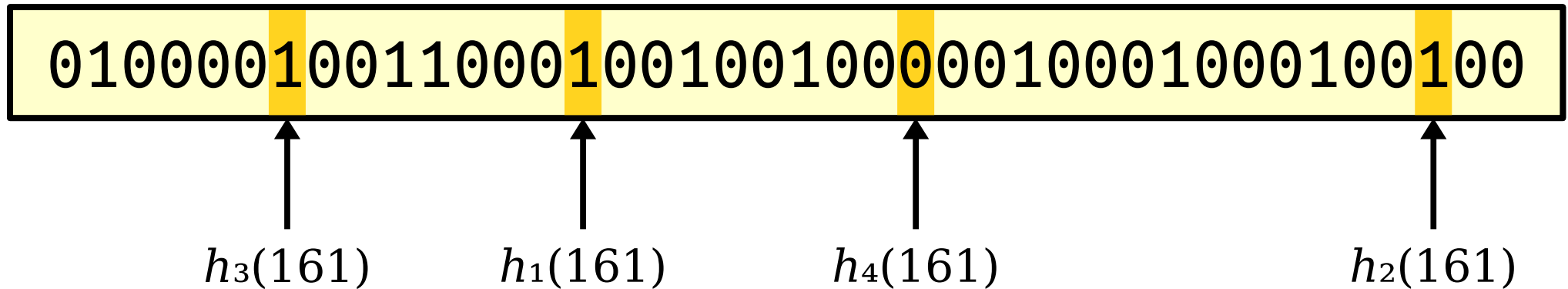


Can we get the same accuracy
while using fewer bits?

create(S): Select k hash functions. Hash each element with all hash functions and set the indicated bits to 1.

query(x): Hash x with all k hash functions.

Return whether all the indicated bits are 1.



Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If c is too low, we'll get too many false positives. If c is too large, we'll use too much memory.

01000010011000100100100000010001000100100

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If c is too low, we'll get too many false positives. If c is too large, we'll use too much memory.

10001000100011001010100001000101

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If c is too low, we'll get too many false positives. If c is too large, we'll use too much memory.

110011001111000111000101

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If c is too low, we'll get too many false positives. If c is too large, we'll use too much memory.

111111

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If c is too low, we'll get too many false positives. If c is too large, we'll use too much memory.

01000010011000100100100000010001000100100

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

01000010011000100100100000010001000100100

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

01010010011010100110100010100010001000110100

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

0101001111101010011010001011001000110111

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

111

Can we get the same accuracy
while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

01000010011000100100100000010001000100100

Can we get the same accuracy while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

00000010010000100000000000000000000010000000100

Can we get the same accuracy
while using fewer bits?

We have two knobs to turn: the coefficient c in the number of bits $cn \lg \varepsilon^{-1}$, and the number of hash functions k .

Intuition: If k is either too low *or* too high, we'll get too many false positives.

Question: How do we tune k , the number of hash functions?

01000010011000100100100000010001000100100

Can we get the same accuracy while using fewer bits?

In what follows, let's have

$$m = cn \lg \epsilon^{-1}$$

be the number of bits in our array.

How do we quantify our error rate?

Question: In what circumstance do we get a false positive?

Answer: Each of the element's bits are set, but the element isn't in the set S .

Question: What is the probability that this happens?

010000**1**0011000**1**00100**1**00000010001000**1**00100

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Focus on a bit at index i .

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Focus on a bit at index i .
Pick some $x \in S$ and hash function h .

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Focus on a bit at index i .
Pick some $x \in S$ and hash function h .

What's the probability that $h(x) \neq i$? (Assume truly random hash functions.)

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Focus on a bit at index i .
Pick some $x \in S$ and hash function h .

What's the probability that $h(x) \neq i$? (Assume truly random hash functions.)

Answer: $1 - 1/m$.

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Focus on a bit at index i .
Pick some $x \in S$ and hash function h .

What's the probability that $h(x) \neq i$? (Assume truly random hash functions.)

Answer: $1 - 1/m$.

What's the probability that, across all n elements and k hash functions, bit i isn't set?

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Focus on a bit at index i .
Pick some $x \in S$ and hash function h .

What's the probability that $h(x) \neq i$? (Assume truly random hash functions.)

Answer: $1 - 1/m$.

What's the probability that, across all n elements and k hash functions, bit i isn't set?

Answer: $(1 - 1/m)^{kn}$.

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Useful fact: $(1 - 1/p)^p \approx e^{-1}$.

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Useful fact: $(1 - 1/p)^p \approx e^{-1}$.

Probability that bit i is unset after inserting n elements:

$$\left(1 - \frac{1}{m}\right)^{kn}$$

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Useful fact: $(1 - 1/p)^p \approx e^{-1}$.

Probability that bit i is unset after inserting n elements:

$$\left(1 - \frac{1}{m}\right)^{kn}$$
$$= \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{kn}{m}}$$

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Useful fact: $(1 - 1/p)^p \approx e^{-1}$.

Probability that bit i is unset after inserting n elements:

$$\begin{aligned} & \left(1 - \frac{1}{m}\right)^{kn} \\ &= \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{kn}{m}} \\ &\approx e^{-\frac{kn}{m}} \end{aligned}$$

How do we quantify our error rate?

Question 1: What is the probability that any particular bit is set?

001101010000101000

Useful fact: $(1 - 1/p)^p \approx e^{-1}$.

Probability that bit i is unset after inserting n elements:

$$\begin{aligned} & \left(1 - \frac{1}{m}\right)^{kn} \\ &= \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{kn}{m}} \\ &\approx e^{-\frac{kn}{m}} \\ &= e^{-k/c \lg \varepsilon^{-1}} \end{aligned}$$

How do we quantify our error rate?

Question 2: What is the probability of a false positive?

00110101000101000

Probability that a fixed bit is 1 after n elements have been added:

$$\approx 1 - e^{-k / c \lg \varepsilon^{-1}}$$

How do we quantify our error rate?

Question 2: What is the probability of a false positive?

00110101000101000

Probability that a fixed bit is 1 after n elements have been added:

$$\approx 1 - e^{-k / c \lg \varepsilon^{-1}}$$

False positive probability is approximately

How do we quantify our error rate?

Question 2: What is the probability of a false positive?

00110101000101000

Probability that a fixed bit is 1 after n elements have been added:

$$\approx 1 - e^{-k / c \lg \varepsilon^{-1}}$$

False positive probability is approximately

How do we quantify our error rate?

Question 2: What is the probability of a false positive?

00110101000101000

Probability that a fixed bit is 1 after n elements have been added:

$$\approx 1 - e^{-k / c \lg \varepsilon^{-1}}$$

False positive probability is approximately

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k$$

How do we quantify our error rate?

Question 2: What is the probability of a false positive?

00110101000101000

This value isn't exactly correct because certain bits being 1 decrease the probability that other bits are 1. With a more advanced analysis we can show that this is very close to the true value.

Probability that a fixed bit is 1 after n elements have been added:

$$\approx 1 - e^{-k / c \lg \varepsilon^{-1}}$$

False positive probability is approximately

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k$$

How do we quantify our error rate?

Question 2: What is the probability of a false positive?


00110101000101000

This value isn't exactly correct because certain bits being 1 decrease the probability that other bits are 1. With a more advanced analysis we can show that this is very close to the true value.

Probability that a fixed bit is 1 after n elements have been added:

$$\approx 1 - e^{-k / c \lg \varepsilon^{-1}}$$

False positive probability is approximately

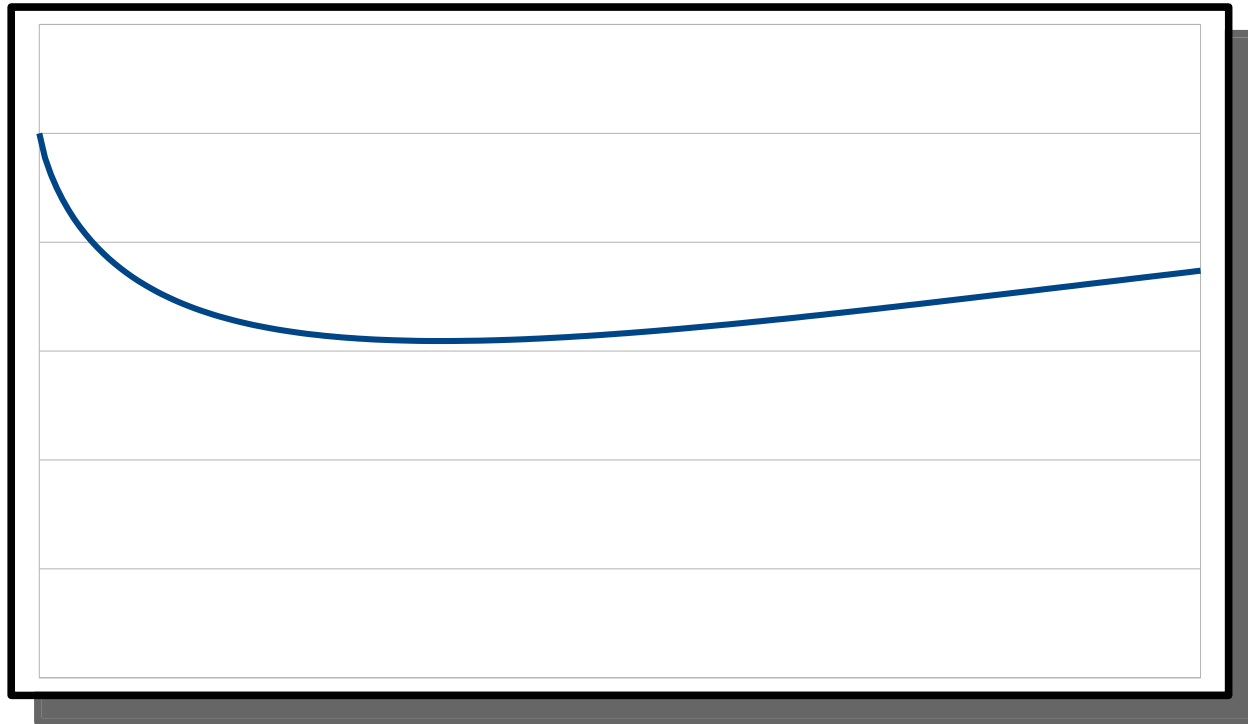

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k$$

Question: What choice of k minimizes this expression?

How do we quantify our error rate?

Goal: Pick k to
minimize

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k.$$



How many hash functions should we use?

Goal: Pick k to minimize

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k.$$

Claim: This expression is minimized when

$$k = c \lg \varepsilon^{-1} \ln 2.$$

You can show this using some symmetry arguments or calculus.

How many hash functions should we use?

Goal: Pick k to minimize

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k.$$

Claim: This expression is minimized when

$$k = c \lg \varepsilon^{-1} \ln 2.$$

You can show this using some symmetry arguments or calculus.

Good exercise: This claim is often repeated and seldom proved. Confirm I am not perpetuating lies.

How many hash functions should we use?

Goal: Pick k to minimize

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k.$$

Claim: This expression is minimized when

$$k = c \lg \varepsilon^{-1} \ln 2.$$

You can show this using some symmetry arguments or calculus.

Good exercise: This claim is often repeated and seldom proved. Confirm I am not perpetuating lies.

Challenge: Give an explanation for this result that is “immediately obvious” from the original expression.

How many hash functions should we use?

The false positive rate is

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k.$$

and we know to pick

$$k = c \lg \varepsilon^{-1} \ln 2.$$

Plugging this value into
the expression gives a
false positive rate of

$$\varepsilon^{c \ln 2}$$

(The derivation, for those of
you who are curious.)

$$(1 - e^{-k / c \lg \varepsilon^{-1}})^k$$

$$= (1 - e^{-c \lg \varepsilon^{-1} \ln 2 / c \lg \varepsilon^{-1}})^k$$

$$= (1 - e^{-\ln 2})^k$$

$$= \left(1 - \frac{1}{2}\right)^k$$

$$= 2^{-k}$$

$$= 2^{-c \lg \varepsilon^{-1} \ln 2}$$

$$= 2^{\lg \varepsilon^{c \ln 2}}$$

$$= \varepsilon^{c \ln 2}$$

Knowing what we know now, how many bits
do we need to get a false positive rate of ε ?

Our false positive rate,
as a function of c , is

$$\varepsilon^{c \ln 2}.$$

Our goal is to get a false
positive rate of ε .

To do so, pick

$$\begin{aligned} c &= 1 / \ln 2 \\ &= \lg e \end{aligned}$$

Which is about 1.44.

Knowing what we know now, how many bits
do we need to get a false positive rate of ε ?

Given a number of elements n and an error rate ε , pick

$$m = cn \lg \varepsilon^{-1}$$

$$k = c \lg \varepsilon^{-1} \ln 2.$$

Optimal c :

$$c = \lg e$$

How did we do overall?

Given a number of elements n and an error rate ε , pick

$$m = n \lg \varepsilon^{-1} \lg e$$

$$k = \lg \varepsilon^{-1}$$

Optimal c :

$$c = \lg e$$

How did we do overall?

Given a number of elements n and an error rate ε , pick

$$m = n \lg \varepsilon^{-1} \lg e$$

$$k = \lg \varepsilon^{-1}$$

Space usage: $\Theta(n \log \varepsilon^{-1})$.
(Within a factor of 1.44 of optimal!)

Query time: $O(\log \varepsilon^{-1})$.

More recent data structures, like the **quotient filter** and **cuckoo filter**, push this coefficient down and are still viable in practice.

We know that it's possible to reduce this coefficient to $1 + o(1)$ (**Pagh, Pagh, and Rao's optimal AMQ structure**), though this is currently only of theoretical interest.

How did we do overall?

To Summarize

- ***Chebyshev's inequality*** is a useful concentration inequality for two-sided errors. It involves bounding the variance of the underlying variable.
- To simplify variance expressions, look for pairwise independence.
- Some random variables (± 1 variables, indicator variables, etc.) simplify nicely when squared.
- Taking the median of many estimators with two-sided error is a good way to amplify success probabilities. The ***Chernoff bound*** is helpful for showing how effective this strategy is.
- Information-theoretic lower bounds can be helpful in determining how much progress can be made with a data structure.

Next Time

- ***Cuckoo Hashing***
 - Hashing with worst-case $O(1)$ lookups.
- ***Random Graph Theory***
 - Properties of random bipartite graphs.