

Hashing and Sketching

Part Two

Outline for Today

- ***Recap from Last Time***
 - Where are we, again?
- ***Count Sketches***
 - A frequency estimator that shows off several key mathematical techniques.
- ***Cardinality Estimators***
 - How many different items have you seen?

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.

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



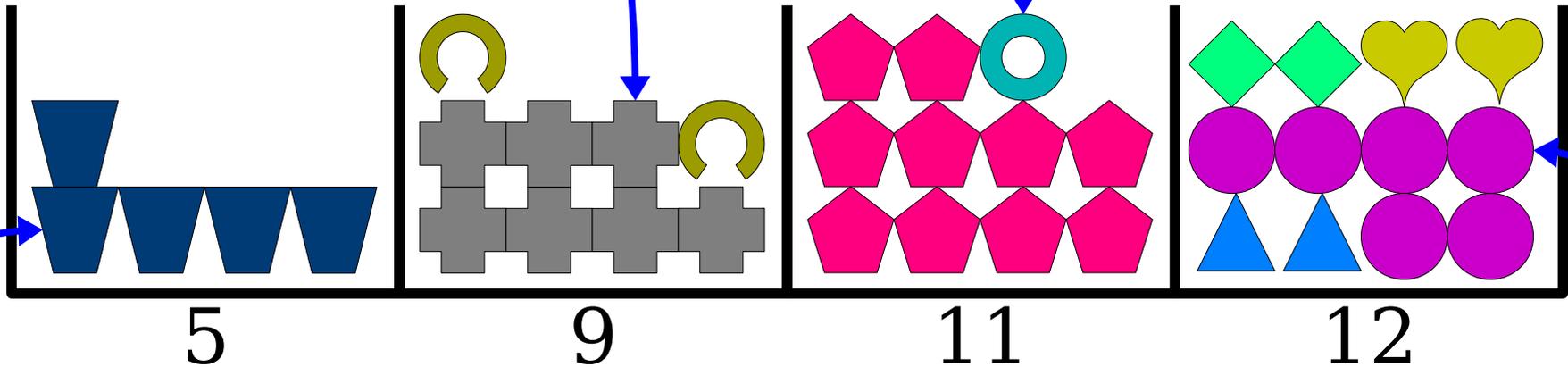
Frequency Estimation

- **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.

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 ().

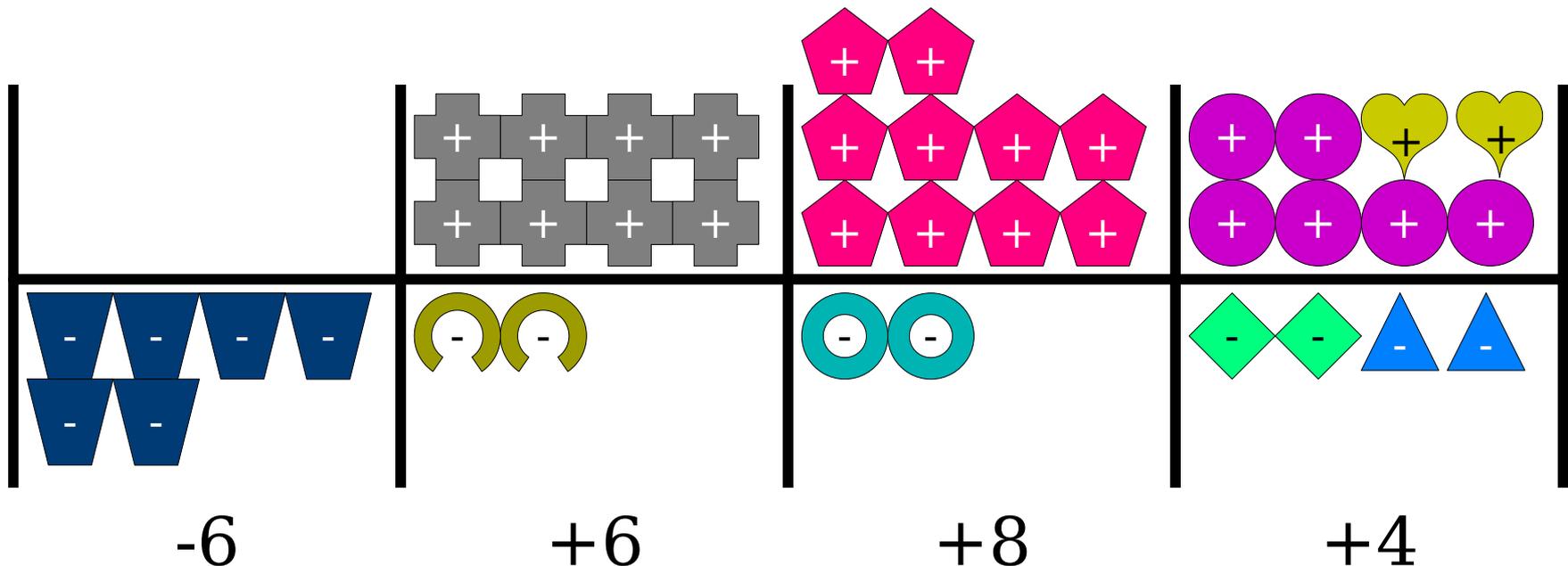


We have a good estimate for , since nothing collides with it.

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

The Setup

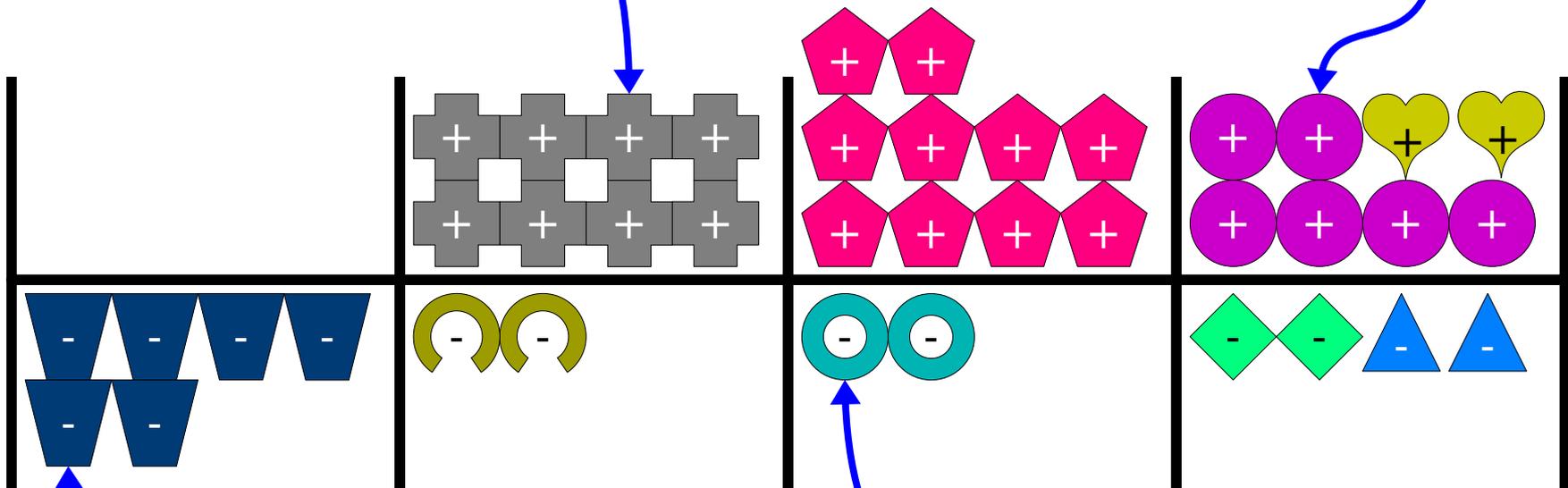
- As before, create an array of counters and assign each item a counter.
- **Key New Step:** For each item x , assign x either $+1$ or -1 .
 - To **increment**(x), go to **count**[$h(x)$] and add ± 1 as appropriate.
 - To **estimate**(x), return **count**[$h(x)$], multiplied by ± 1 as appropriate.



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 are mostly offset.

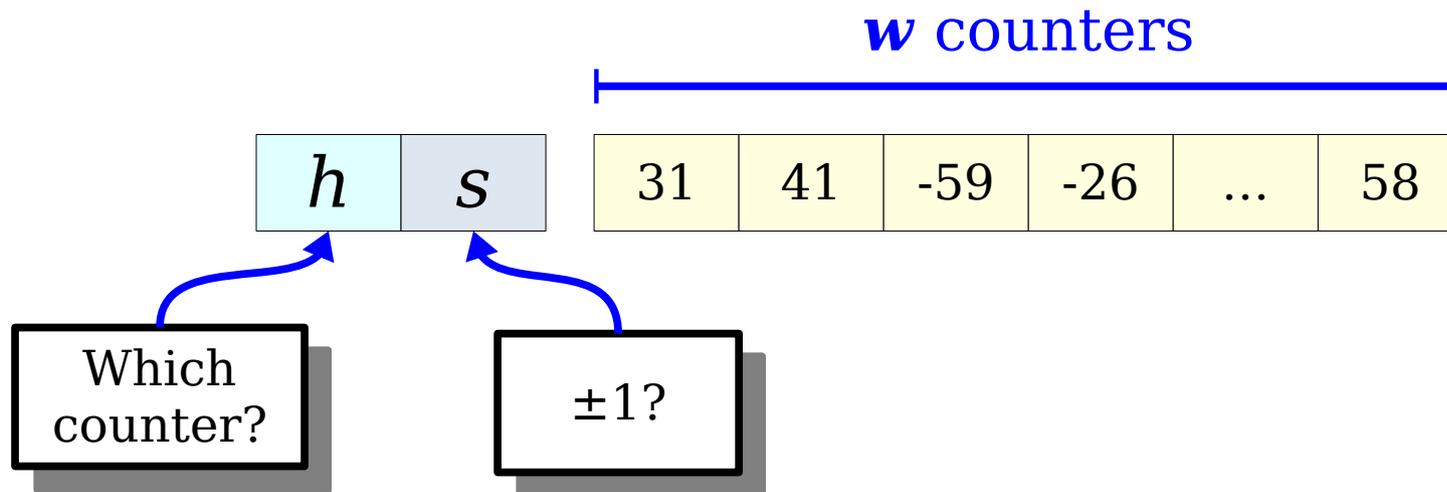


We have a good estimate for , since nothing collides with it.

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

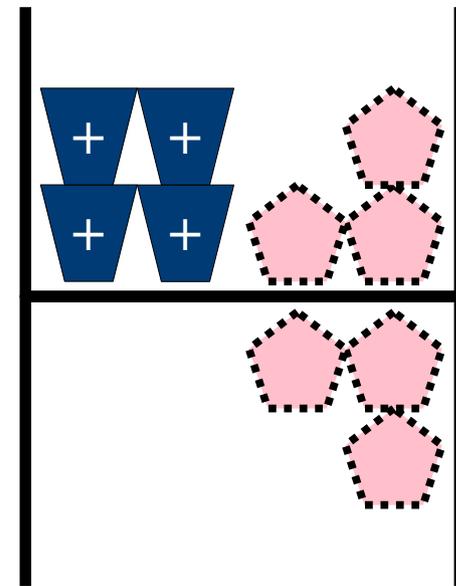
Formalizing This

- 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 $s \in \mathcal{S}$ 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 $s(x) \cdot \text{count}[h(x)]$.



The Expectation, Intuitively

- Focus on any element x_i whose frequency we're estimating.
- Think about any element that collides with us.
- With 50% probability, it *increases* our estimate.
- With 50% probability, it *decreases* our estimate.
- **Intuition:** The expected value weights both options equally, so our estimator will be unbiased.



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 introduce random variables indicating whether x_i and x_j collided with one another:

$$\mathbb{1}_{h(x_i)=h(x_j)} = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{if } h(x_i) \neq h(x_j) \end{cases}$$

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

$$\hat{\mathbf{a}}_i = \sum_j \mathbf{a}_j s(x_i) s(x_j) \mathbb{1}_{h(x_i)=h(x_j)}$$

How much the collision hurts us...

... if there even is a collision at all.

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{a}}_i] &= \mathbb{E}\left[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j s(x_i) s(x_j) \mathbb{1}_{h(x_i)=h(x_j)}\right] \\
&= \mathbb{E}[\mathbf{a}_i] + \mathbb{E}\left[\sum_{j \neq i} \mathbf{a}_j s(x_i) s(x_j) \mathbb{1}_{h(x_i)=h(x_j)}\right] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}\left[\mathbf{a}_j s(x_i) s(x_j) \mathbb{1}_{h(x_i)=h(x_j)}\right] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(x_i) s(x_j)] \mathbb{E}[\mathbf{a}_j \mathbb{1}_{h(x_i)=h(x_j)}] \\
&= \mathbf{a}_i + \sum_{j \neq i} \mathbb{E}[s(x_i)] \mathbb{E}[s(x_j)] \mathbb{E}[\mathbf{a}_j \mathbb{1}_{h(x_i)=h(x_j)}] \\
&= \mathbf{a}_i + \sum_{j \neq i} 0 \\
&= \mathbf{a}_i
\end{aligned}$$

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

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

$s(x_i)$ is uniform over $\{-1, +1\}$

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

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

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

- 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.

$$\begin{aligned}\text{Var}[\hat{\mathbf{a}}_i] &= \text{Var}\left[\mathbf{a}_i + \sum_{j \neq i} \mathbf{a}_j \mathbf{s}(\mathbf{x}_i) \mathbf{s}(\mathbf{x}_j) \mathbb{1}_{h(\mathbf{x}_i)=h(\mathbf{x}_j)}\right] \\ &= \text{Var}\left[\sum_{j \neq i} \mathbf{a}_j \mathbf{s}(\mathbf{x}_i) \mathbf{s}(\mathbf{x}_j) \mathbb{1}_{h(\mathbf{x}_i)=h(\mathbf{x}_j)}\right] \\ &= \sum_{j \neq i} \text{Var}\left[\mathbf{a}_j \mathbf{s}(\mathbf{x}_i) \mathbf{s}(\mathbf{x}_j) \mathbb{1}_{h(\mathbf{x}_i)=h(\mathbf{x}_j)}\right]\end{aligned}$$

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 pairwise uncorrelated.
(*Prove this!*)

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

The techniques here transfer well to other problems!

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?

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

This is the square of the magnitude of the vector.

The magnitude of a vector is often denoted $\|\mathbf{a}\|$.

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

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

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

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

Where We Stand

- We know that

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

- With the count-min sketch, we bounded the probability that we overestimated by more than $\varepsilon\|\mathbf{a}\|_1$.
- Since the variance is related to $\|\mathbf{a}\|_2$, for the count sketch we'll bound the probability that we are more than $\varepsilon\|\mathbf{a}\|_2$ from our estimate:

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

$$\begin{aligned} & \Pr[|\hat{\mathbf{a}}_i - \mathbf{a}_i| > \varepsilon \|\mathbf{a}\|_2] \\ & \leq \frac{\text{Var}[\hat{\mathbf{a}}_i]}{(\varepsilon \|\mathbf{a}\|_2)^2} \\ & \leq \frac{\|\mathbf{a}\|_2^2}{w} \cdot \frac{1}{(\varepsilon \|\mathbf{a}\|_2)^2} \\ & = \frac{1}{w \varepsilon^2} \end{aligned}$$

Taking Stock

- We've just shown that

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

- This means that we need $w = \Theta(\varepsilon^{-2})$ in order to get a strong bound.
 - Compare with the count sketch, where $w = \Theta(\varepsilon^{-1})$.
- **Idea:** Set $w = 4 \cdot \varepsilon^{-2}$.
 - Why 4? Because I peeked ahead.

Running in Parallel

- Imagine we run d copies of this estimator and call *estimate*(x) on each of our estimators and get back these estimates.
- We need to give back a single number.
- **Question:** How should we aggregate these numbers into a single estimate?

Answer at
<https://cs166.stanford.edu/pollev>

Estimator 1:
137

Estimator 2:
271

Estimator 3:
166

Estimator 4:
103

Estimator 5:
261

Running in Parallel

- Unlike last time, we have a two-sided error, so taking the minimum would be a Very Bad Idea.
- Two reasonable options come to mind:
 - Take the **mean** of the estimates.
 - Take the **median** of the estimates.
- **Question:** Which should we pick?

Estimator 1:
137

Estimator 2:
271

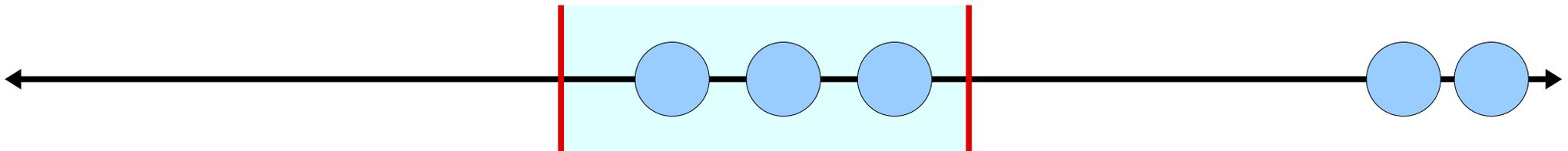
Estimator 3:
166

Estimator 4:
103

Estimator 5:
261

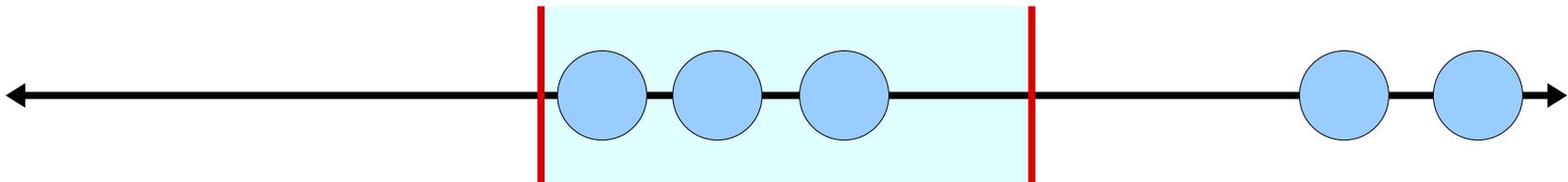
Working With Means

- **Claim:** Taking the mean of multiple estimators does increase our probability of being close to the expected value, but not very quickly.
- **Intuition:** Not all outliers are created equal, and outliers far from the target range skew the estimate.
- **The Math:** Averaging d copies of an estimator decreases the variance by a factor of d . (Prove this!) By Chebyshev, that decreases the probability of getting a bad answer by a factor of d . We'd like something that decays exponentially in d .



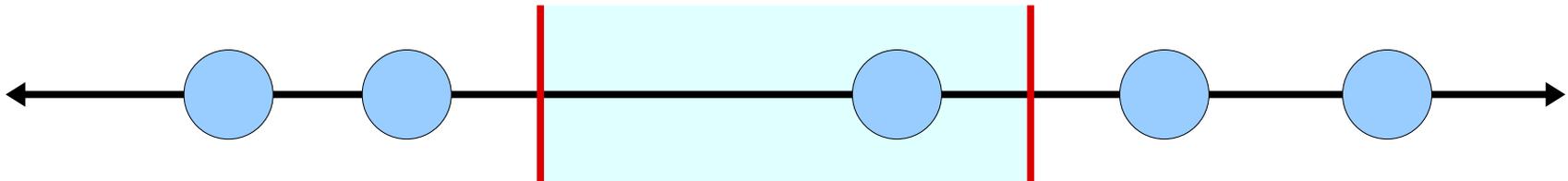
Working With Medians

- **Claim:** If we output the median estimate given by the data structures, we have high probability of giving an acceptably close 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 Medians

- 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

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

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

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

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

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

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

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

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

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

$$w = \lceil 4 \cdot \epsilon^{-2} \rceil$$

h_1	s_1
h_2	s_2
h_3	s_3
...	...
h_d	s_d

31	41	-59	-26	...	58
27	-18	28	-18	...	-45
16	-18	-3	39	...	-75
...					
69	-31	47	-18	...	59

$$d = \lceil 8 \ln \frac{1}{\epsilon} \rceil$$

Sampled uniformly and independently from 2-independent families of hash functions

The Count Sketch

$$w = \lceil 4 \cdot \epsilon^{-2} \rceil$$

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

Diagram annotations:
- A blue horizontal bracket above the second table spans the width of the sketch, labeled $w = \lceil 4 \cdot \epsilon^{-2} \rceil$.
- A blue vertical bracket to the right of the second table spans the height of the sketch, labeled $d = \lceil 8 \ln \frac{1}{\epsilon} \rceil$.

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

The Count Sketch

$$w = \lceil 4 \cdot \epsilon^{-2} \rceil$$

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

$d = \lceil 8 \ln \frac{1}{\epsilon} \rceil$

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

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

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

Accuracy: within $\varepsilon \|\mathbf{a}\|_1$.

Count Sketch

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

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

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

Accuracy: within $\varepsilon \|\mathbf{a}\|_2$

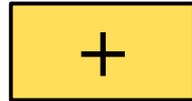
Major Ideas Here

- Concentration inequalities are useful tools for showing the right thing probably happens.
 - For one-sided errors, try Markov's inequality.
 - For two-sided errors, try Chebyshev's inequality.
 - To bound the probability that lots of things all go wrong, use Chernoff bounds.
 - For more on different mathematical tools like these, check out [*this blog post by Scott Aaronson*](#).
- Modest success probability can be amplified by running things in parallel.
 - For one-sided errors, try using the min or max.
 - For two-sided errors, try using the median.
- We can estimate quantities using significantly less space than storing those quantities exactly if we're okay with approximate answers.

Cardinality Estimation

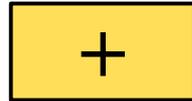
Cardinality Estimation

- A **cardinality estimator** is a data structure supporting the following operations:
 - **see**(x), which records that x has been seen, and
 - **estimate**(), which returns an estimate of the number of **distinct** values we've seen.
- In other words, they estimate the cardinality of the set of all items that have been seen.
- These data structures are widely deployed in practice.
 - Databases use them to select which of many different algorithms to run, based on the number of items to process.
 - Websites use them to estimate how many different people have visited the site in a given time window.



Cardinality Estimation

- As with frequency estimation, we can solve the cardinality estimation problem exactly using hash tables or binary search trees using $\Omega(n)$ space.
- To be useful in large-scale data applications, cardinality estimators need to use *significantly* less space than this.
- **Question:** How low can we go?



Flipping Coins

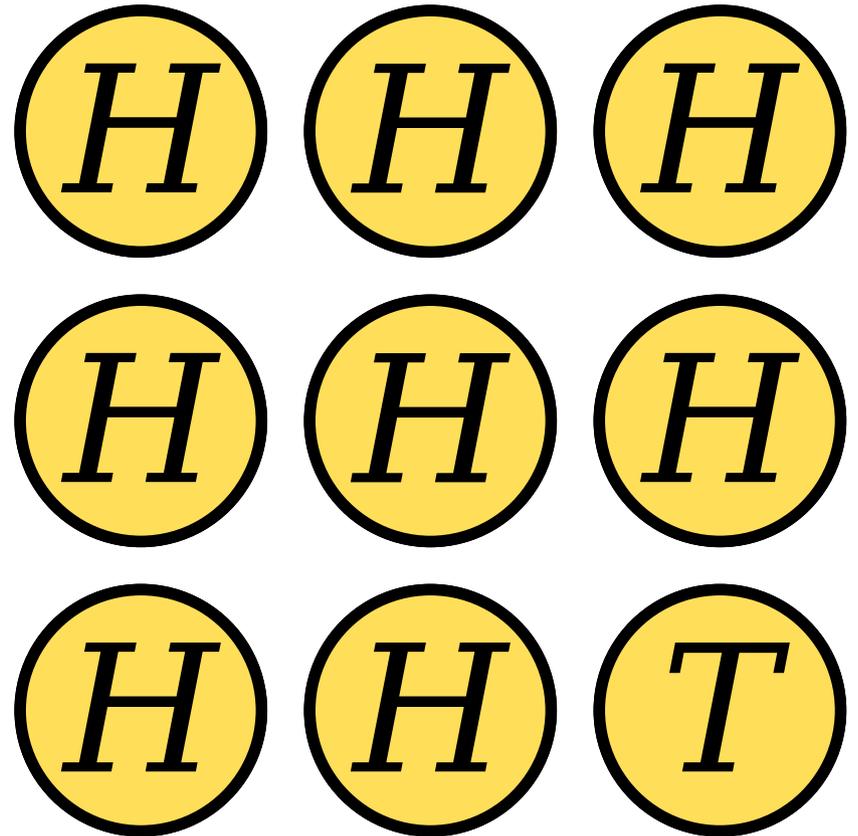
Flipping Coins

- Here's a game: I'm going to flip a coin until I get tails. My score is the number of coins that I flip.
- The probability of flipping k or more consecutive heads is 2^{-k} , so it's pretty unlikely that I'm going to flip lots of heads in a row.



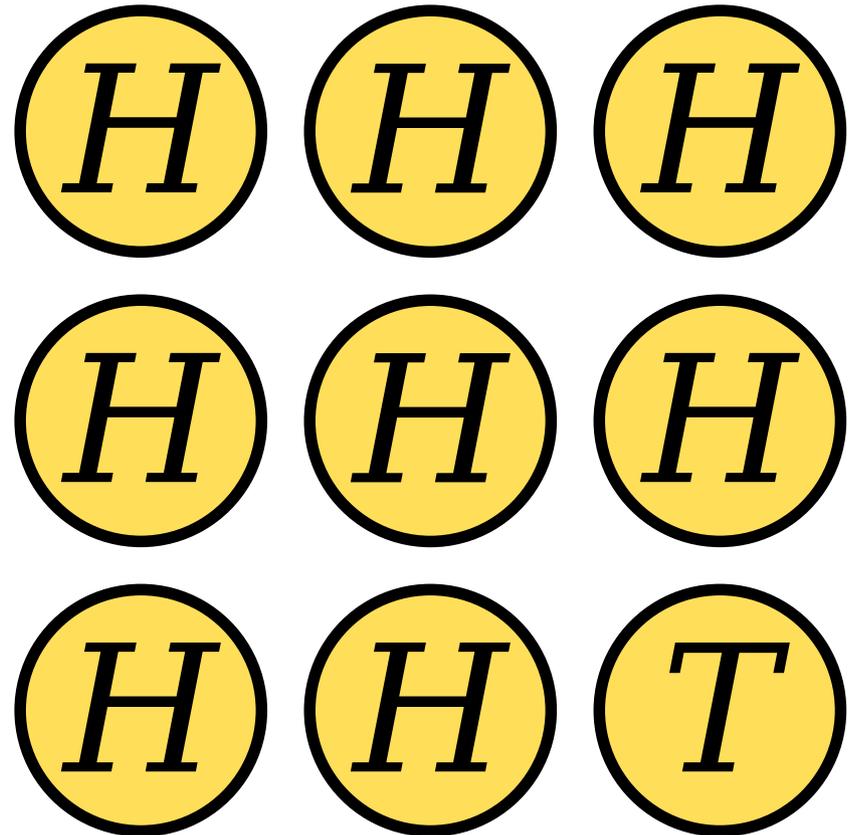
Flipping Coins

- Suppose I show you the following clip of me playing this game.
- Which is more likely?
 - I played the game once and got really lucky.
 - I played the game 512 times and showed you my best run.
- Probability you see this after one game: $1/512$.
- Probability this is the best you see after 512 games: $\approx 23\%$.



Flipping Coins

- **Intuition:** Play this game multiple times and track the maximum score.
- If our maximum score is S , estimate that we played 2^S times.
- **Question:** How good of an estimate is this?

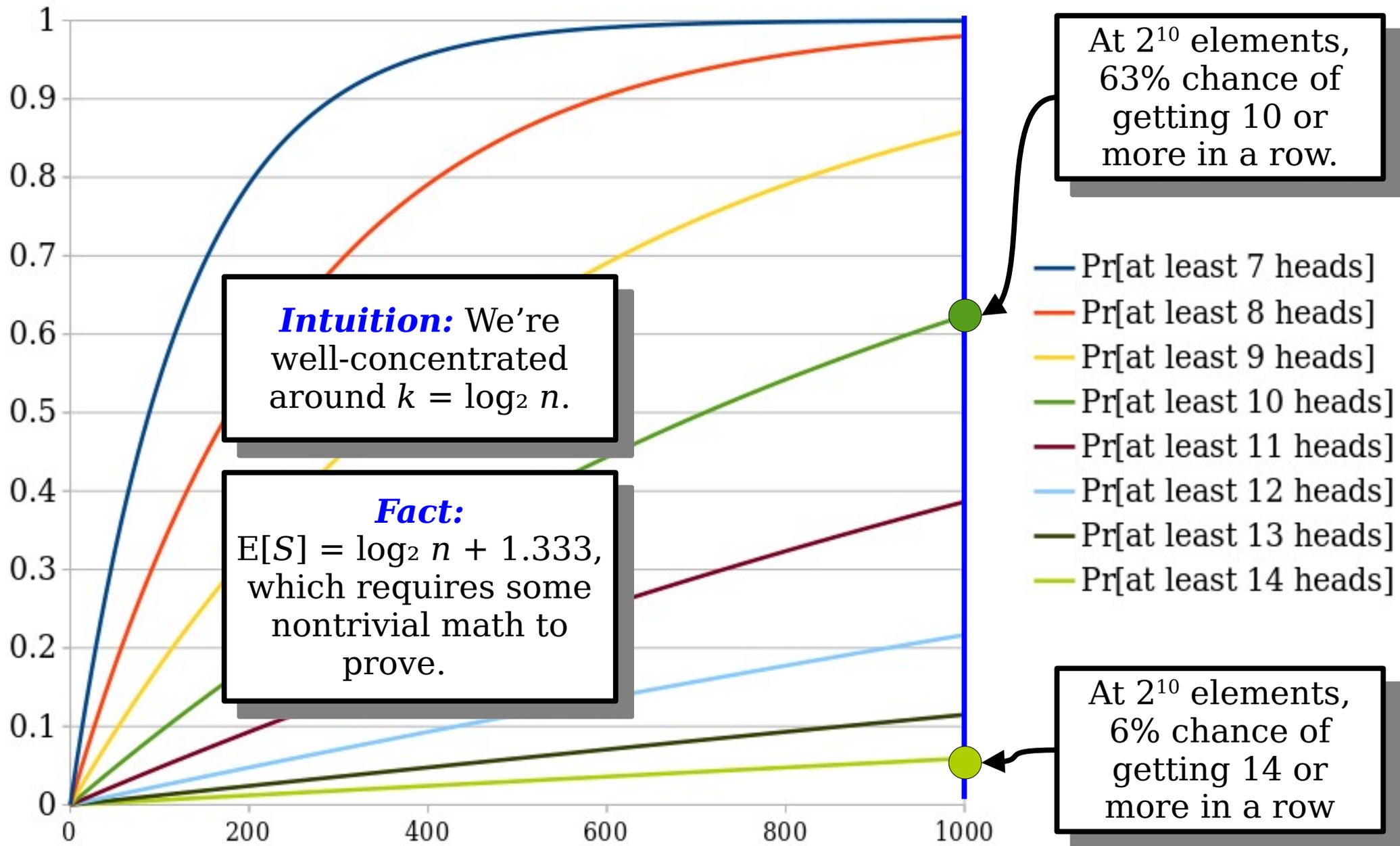


Flipping Coins

- Suppose we play this game n times. What's the probability we see at least k consecutive heads at least once?

$$\begin{aligned} & \Pr[\text{see at least } k \text{ heads in } n \text{ games}] \\ &= 1 - \Pr[\text{never see } k \text{ heads in } n \text{ games}] \\ &= 1 - \Pr[\text{never see } k \text{ heads in one game}]^n \\ &= 1 - (1 - 2^{-k})^n \end{aligned}$$

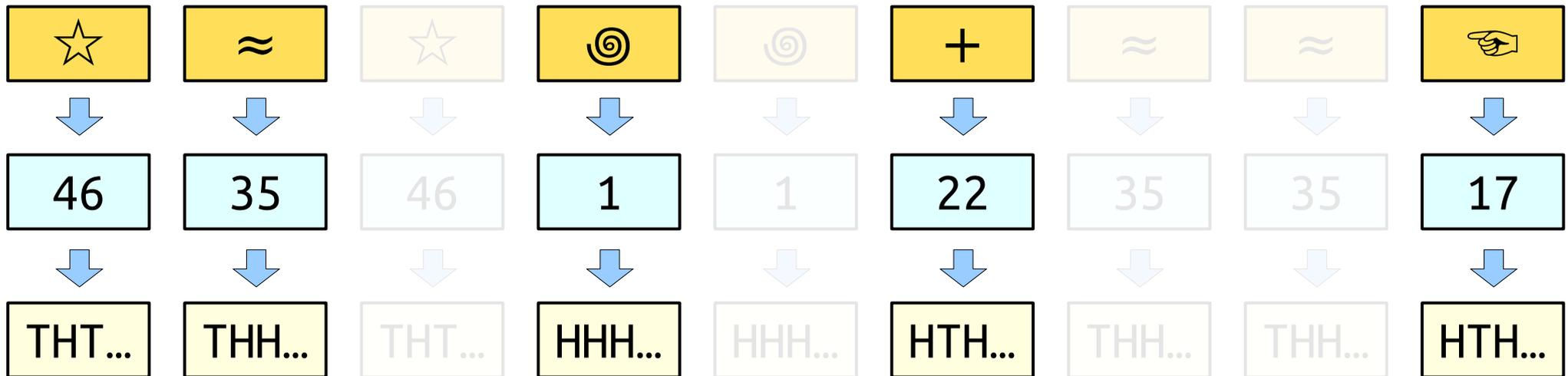
- What does this function look like?



Play this game n times. What is the probability that we see at least k consecutive heads?

From Coins to Cardinality

- Ultimately, we're interested in building a cardinality estimator. How does this help us?
- **Idea:** Hash each item in the data stream, and use each hash as the random source for the coin-flipping game.
- Duplicate items give duplicate hashes, which provide duplicate games, which function as if they never happened.
- If we track the highest score across all these games, we can use that to estimate how many games we played, which is equal to the number of distinct elements we saw.



From Coins to Cardinality

- We need some way of going from hash codes to sequences of coin tosses.
- **Idea:** Treat the hash as a sequence of bits. 0 means heads, 1 means tails.
- Then, count how many appear at or after the last 1 bit in the number.



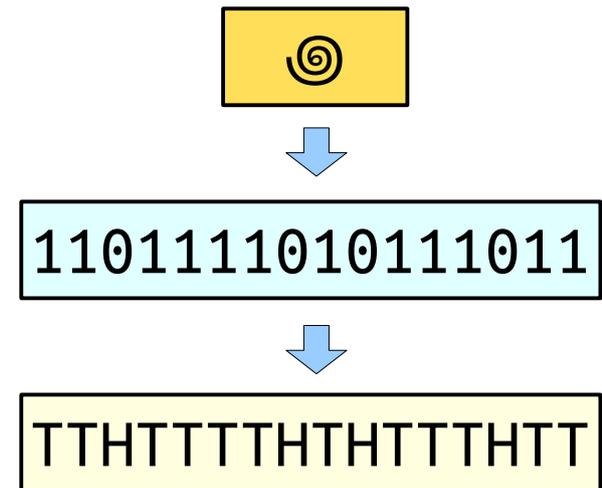
10011111110111101011101100100000



THHTTTTTTHTTTTHTHTTTHTTHTHHHHH

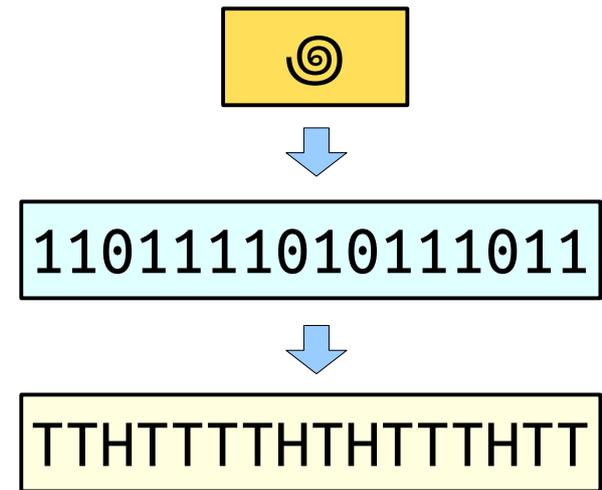
A Simple Estimator

- Keep track of a value S , initially zero, that records how many coins are flipped in a game.
- To **see** an item:
 - Compute a hash code for that item.
 - Compute the index of the lowest 1 bit, plus 1 (the number of coins tossed in our simulated game.)
 - Update S if this is a new record.
- To **estimate** the number of distinct elements:
 - Return 2^S .



A Simple Estimator

- How much space does this single estimator need?
- Assume we have an upper bound U on the maximum cardinality. Our hashes only need $\Theta(\log U)$ bits.
- Bits required to write down the position of a bit in that hash: $\Theta(\log \log U)$.
- That is an *absolutely tiny* amount of space compared to storing the elements!



Improving the Estimator

- The current estimator has a few weaknesses.
 - It always outputs a size that's a power of two, so we're likely to be off by a full binary order of magnitude.
 - It tends to skew high, since a single unexpected run of heads pushes the whole total up.
- But we have already seen some techniques for improving estimators:
 - Run lots of copies in parallel to reduce the likelihood of any one of them being bad.
 - Use some creative strategy to combine those individual estimates into one really good one.
- And in fact, folks have done just that.

HyperLogLog

- The **HyperLogLog** estimator uses many independent copies of this estimator to produce a very high-quality estimate.
 - Run m copies of the estimator, using a hash function to distribute items to estimators, so that each copy gets roughly a $1/m$ fraction of the items.
 - Compute the *harmonic mean* of the estimates to mitigate outliers while smoothing between powers of two.
 - Multiply in a debiasing term to mitigate the skew from both the original estimates and the harmonic mean.
- This estimator is used extensively in practice; with under 1kB of memory, it can estimate cardinalities for any real-world data stream to about 3% accuracy.
- It's widely used in database systems, and many open-source implementations are available.

HyperLogLog

- The analysis of HyperLogLog from the original paper is exceedingly difficult, and I haven't been able to follow along with all the details.
- Hopefully, this intuitive explanation of how it works is enough for you.
- ***(Probably?) Open problem:*** Find a significantly simpler and cleaner rigorous analysis of HyperLogLog than the original.

Major Ideas We've Seen

- You can build a great estimator by running lots of weak estimators in parallel and aggregating the results.
- Indicator variables and linearity of expectation are powerful tools when analyzing sketches.
- Markov's and Chebyshev's inequalities are useful for bounding probabilities involving hashing.
- The Chernoff bound is a great tool for showing it's unlikely for lots of things to go wrong.

Next Time

- ***Balanced Trees***
 - How to keep binary search trees balanced.
- ***Red/Black Trees***
 - The classic balanced tree.
- ***B-Trees***
 - Another classic workhorse.