# Hashing and Sketching

Part One

# Outline for Today

- ***Hash Functions***
  - Understanding our basic building blocks.
- ***Frequency Estimation***
  - Estimating how many times we've seen something.
- ***Concentration Inequalities***
  - "Correct on expectation" versus "correct with high probability."
- ***Probability Amplification***
  - Increasing our confidence in our answers.

# Preliminaries: *Hash Functions*

# Hashing in Practice

- Hash functions are used extensively in programming and software engineering:

  - They make hash tables possible: think C++ `std::hash`, Python's `__hash__`, or Java's `Object.hashCode()`.

  - They're used in cryptography: SHA-256, HMAC, etc.

- ***Question:*** When we're in Theoryland, what do we mean when we say "hash function?"

# Hashing in Theoryland

- In Theoryland, a hash function is a function from some domain called the **universe** (typically denoted $\mathcal{U}$) to some codomain.

- The codomain is usually a set of the form

$$[m] = \{0, 1, 2, 3, \ldots, m - 1\}$$

$$h : \mathcal{U} \to [m]$$

# Hashing in Theoryland

- ***Intuition:*** No matter how clever you are with designing a specific hash function, that hash function isn't random, and so there will be pathological inputs.

  - You can formalize this with the pigeonhole principle.

- ***Idea:*** Rather than finding the One True Hash Function, we'll assume we have a collection of hash functions to pick from, and we'll choose which one to use randomly.

# Families of Hash Functions

- A ***family*** of hash functions is a set $\mathcal{H}$ of hash functions with the same domain and codomain.

- We can then introduce randomness into our data structures by sampling a random hash function from $\mathcal{H}$.

- ***Key Point:*** The randomness in our data structures almost always derives from the random choice of hash functions, not from the data.
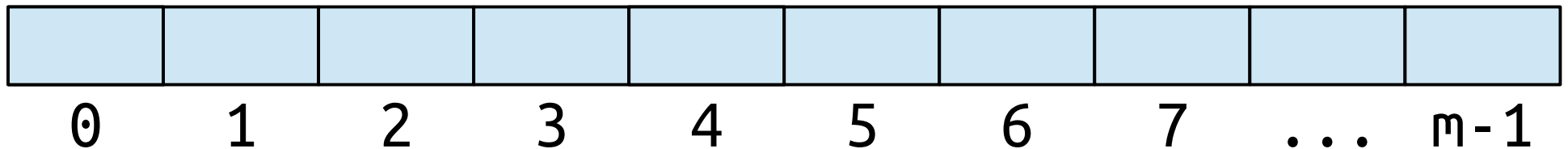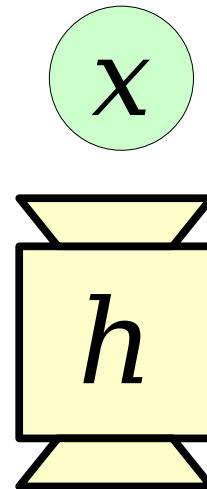
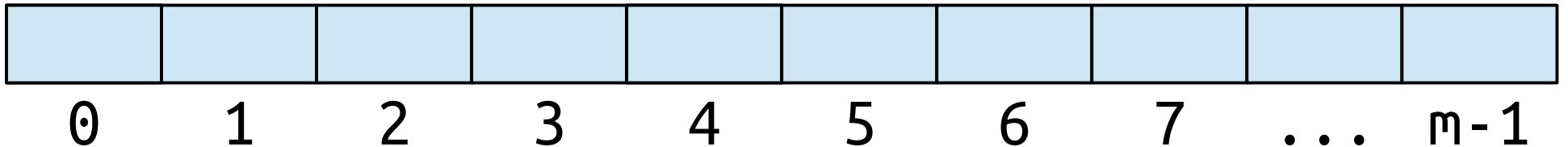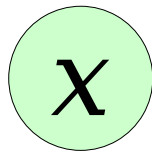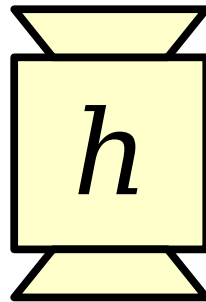**Data is adversarial.**
**Hash function selection is random.**

- ***Question:*** What makes a family of hash functions $\mathcal{H}$ a "good" family of hash functions?

**Goal:** If we pick $h \in \mathscr{H}$ uniformly at random, then $h$ should distribute elements uniformly randomly.

$x$

$h$

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |

$h$

$x$

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |

**Goal:** If we pick $h \in \mathcal{H}$ uniformly at random, then $h$ should distribute elements uniformly randomly.

$y$

$h$

$x$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |

$h$

$x$

$y$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | . . . | m-1 |

$z$

$h$

$x$

$y$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |

**Goal:** If we pick $h \in \mathcal{H}$ uniformly at random, then $h$ should distribute elements uniformly randomly.

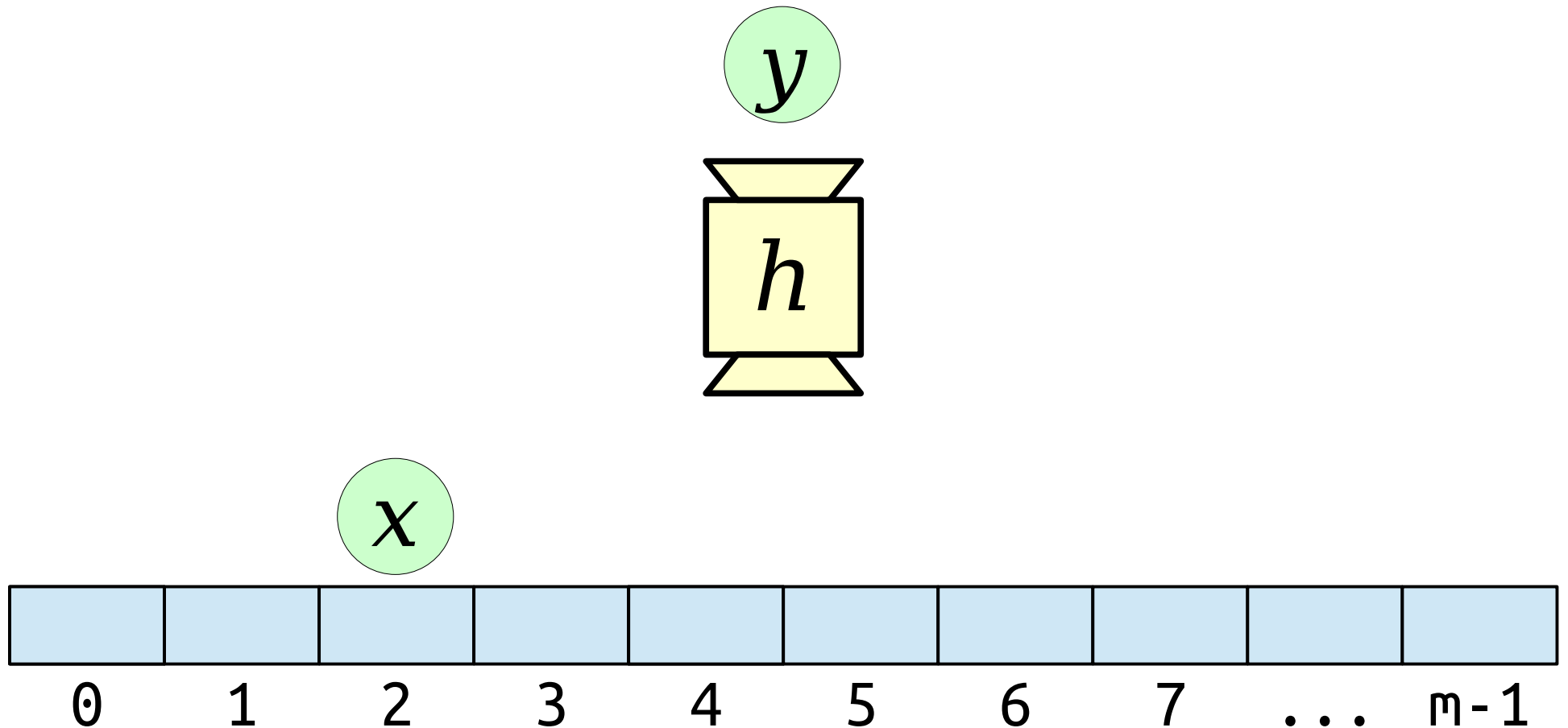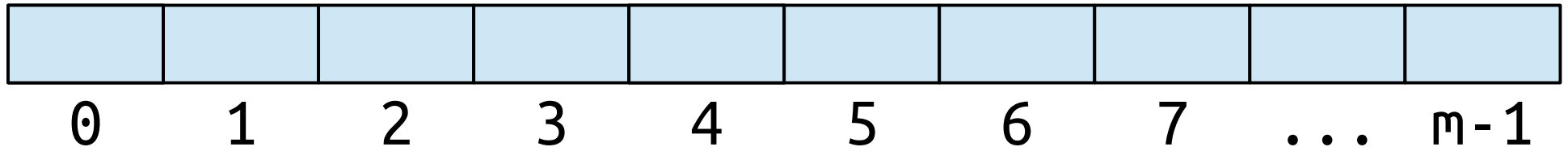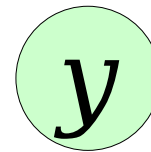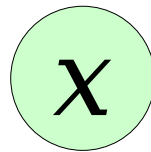**Problem:** A hash function that distributes $n$ elements uniformly at random over $[m]$ requires $\Omega(n \log m)$ space in the worst case.

**Question:** Do we actually need true randomness? Or can we get away with something weaker?

$h$

$z$

$x$

$y$

0   1   2   3   4   5   6   7   ...   m-1

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

x

$\mathcal{H}$

x

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |

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

$x$

$\mathcal{H}$

$x$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |

**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 $[m]$.

Find an "obviously bad" family of hash functions that satisfies the distribution property.

Formulate a hypothesis, but **don't post anything in chat just yet**.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |
|---|---|---|---|---|---|---|---|-----|-----|

**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 $[m]$.

Find an "obviously bad" family of hash functions that satisfies the distribution property.

Now, *private chat me your best guess*. Not sure? Just answer "??".

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |
|---|---|---|---|---|---|---|---|-----|-----|

**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 $[m]$.

**Problem:** This rule doesn't guarantee that elements are spread out.

$w$

$z$

$y$

$x$

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | |

0  1  2  3  4  5  6  7  ... m-1

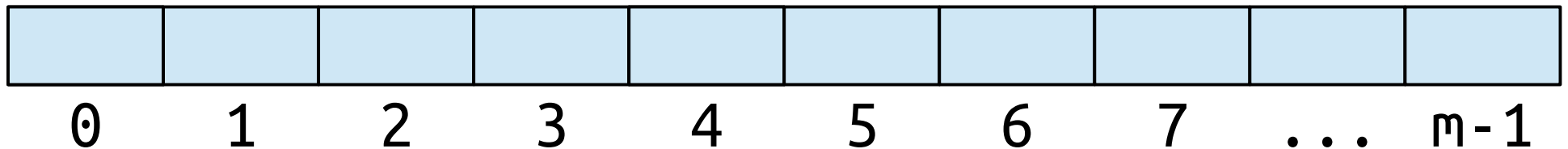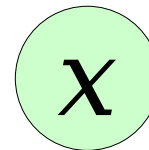**Distribution Property:** Each element should have an equal probability of being placed in each slot.

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

**Problem:** This rule doesn't guarantee that elements are spread out.

$w$

$z$

$y$

$x$

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

0   1   2   3   4   5   6   7   ...   m-1

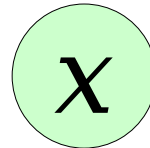**Distribution Property:** Each element should have an equal probability of being placed in each slot.

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

**Problem:** This rule doesn't guarantee that elements are spread out.



0     1     2     3     4     5     6     7     . . .   m-1

**Distribution Property:**
Each element should have an equal probability of being placed in each slot.

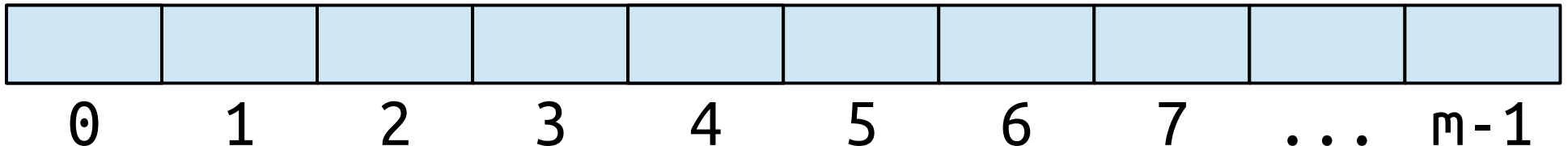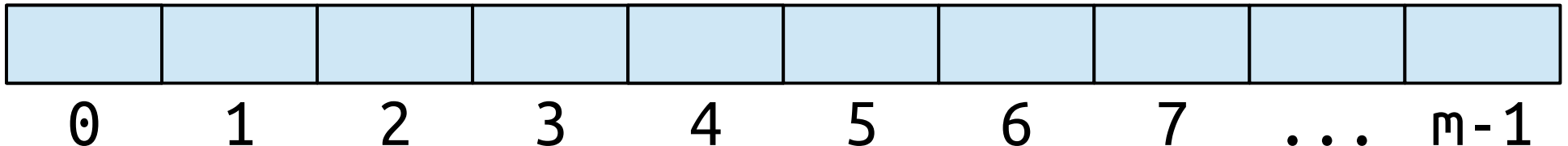For any $x \in \mathscr{U}$ and random $h \in \mathscr{H}$, the value of $h(x)$ is uniform over $[m]$.

**Independence Property:**
Where one element is placed shouldn't impact where a second goes.

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

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

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | ... | m-1 |
|---|---|---|---|---|---|---|---|-----|-----|

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$y$

$x$

| | | | | |
|---|---|---|---|---|

0    1    2    ...   m-1

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$y$

$x$

| | | | | |
|---|---|---|---|---|

0    1    2    ...    m-1

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

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

*Intuition:*
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

*Question:* Where did these elements collide with one another?



0   1   2   ... m-1

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

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

**Intuition:** 2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

**Question:** Where did these elements collide with one another?

$y$

$x$

| | | | | |
|---|---|---|---|---|

0    1    2    . . .   m-1

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

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

***Intuition:*** 2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

***Question:*** Where did these elements collide with one another?

$y$

$x$

0   1   2   ...   m-1

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

*Question:* Where did these elements collide with one another?

$y$

$x$

| | | | | |
|---|---|---|---|---|

0    1    2    ...  m-1

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

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

*Intuition:* 2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

*Question:* Where did these elements collide with one another?



$y$

$x$

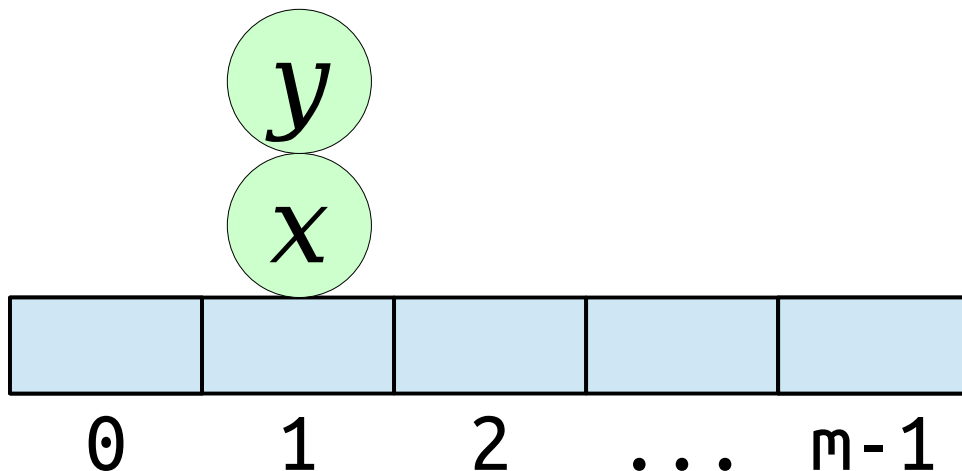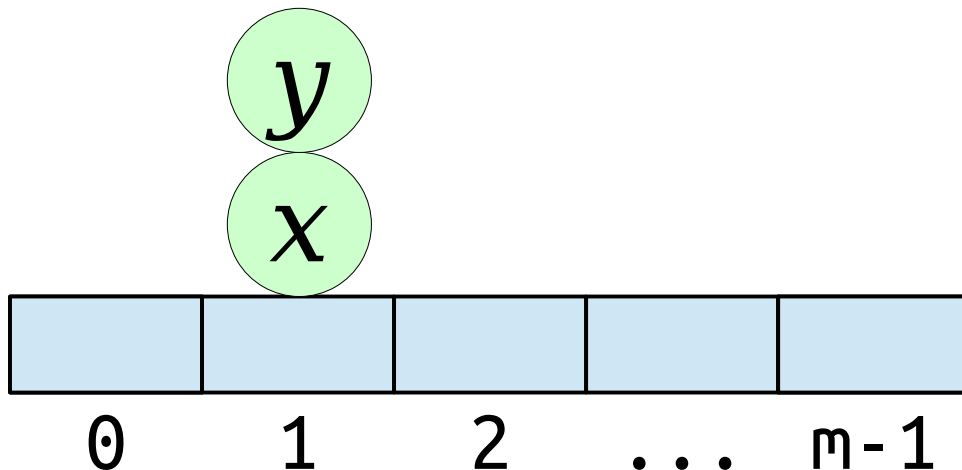| | | | | |
|---|---|---|---|---|

0    1    2    ...    m-1

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

y
x

0    1    2    ...    m-1

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

y

x

| | | | | |
|---|---|---|---|---|

0  1  2  ...  m-1

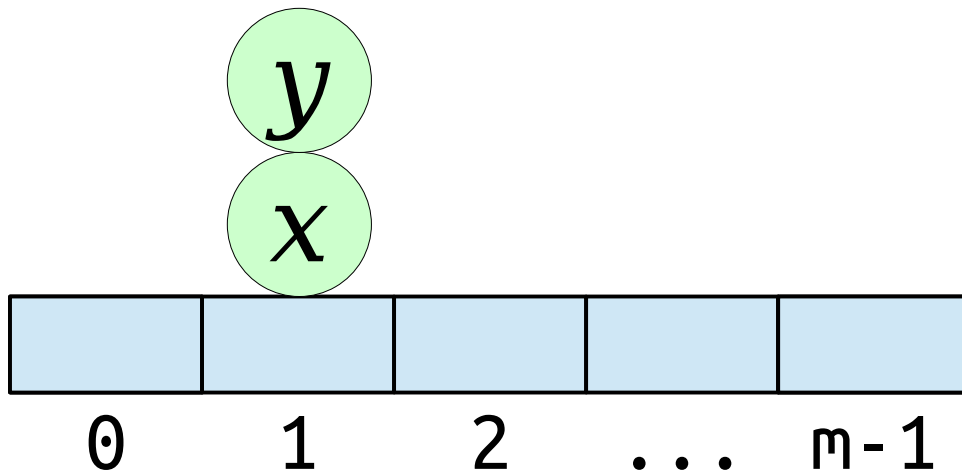For any $x \in \mathcal{U}$ and random $h \in \mathcal{H}$, the value of $h(x)$ is uniform over $[m]$.

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

$y$
$x$

0    1    2    ...  m-1

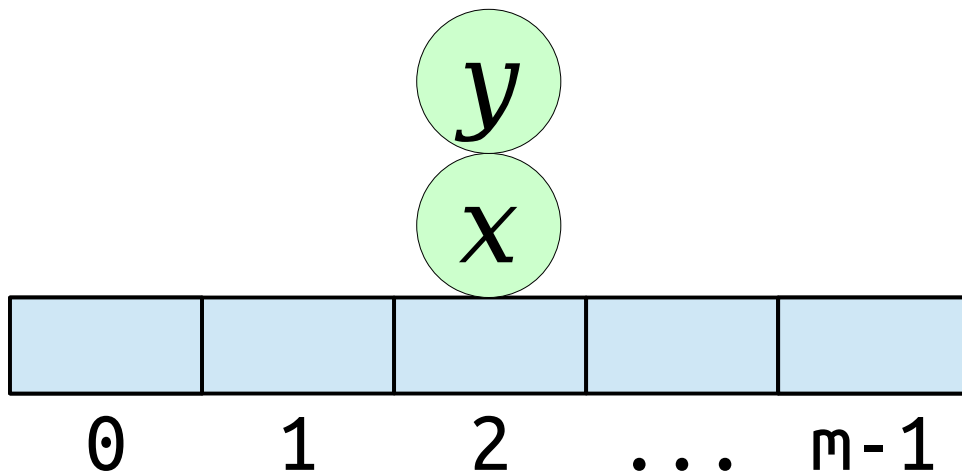For any $x \in \mathcal{U}$ and random $h \in \mathcal{H}$, the value of $h(x)$ is uniform over $[m]$.

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

*Intuition:*
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

$y$
$x$

| | | | | |
|---|---|---|---|---|

0   1   2  ...  m-1

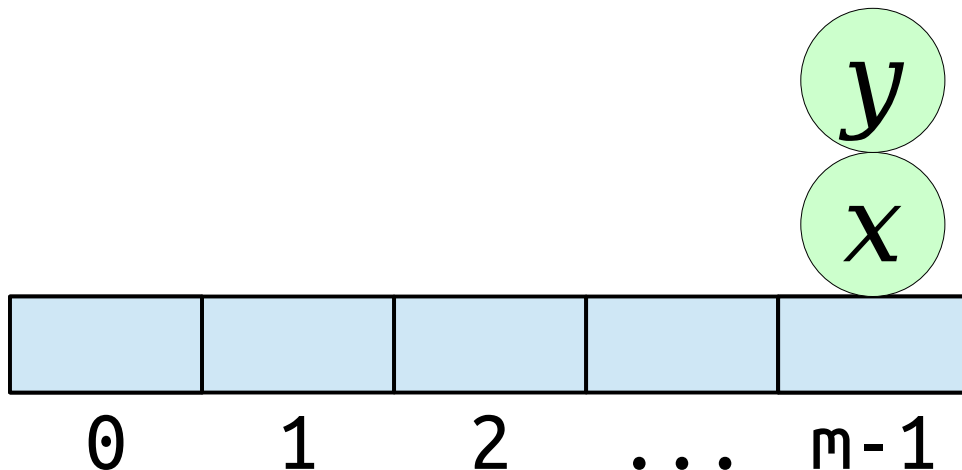For any $x \in \mathcal{U}$ and random $h \in \mathcal{H}$, the value of $h(x)$ is uniform over $[m]$.

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

***Intuition:***
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

$y$

$x$

0    1    2    ...    m-1

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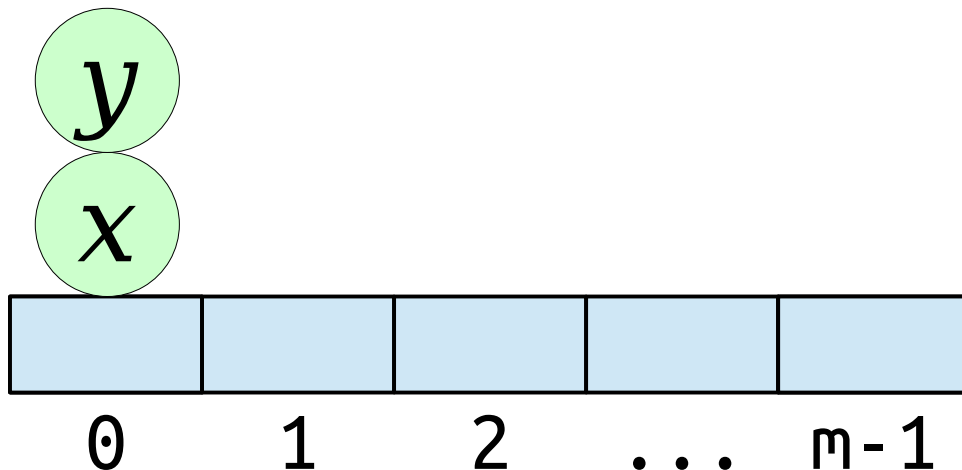
For any distinct $x, y \in \mathcal{U}$ and random $h \in \mathcal{H}$, $h(x)$ and $h(y)$ are independent random variables.

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

$$= \sum_{i=0}^{m-1} \frac{1}{m^2}$$



0    1    2    ...    m-1

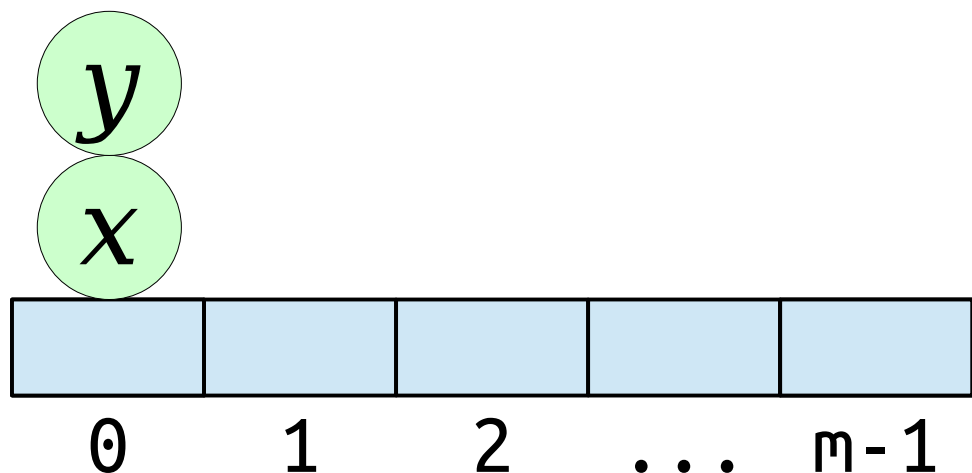For any $x \in \mathcal{U}$ and random $h \in \mathcal{H}$, the value of $h(x)$ is uniform over $[m]$.

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

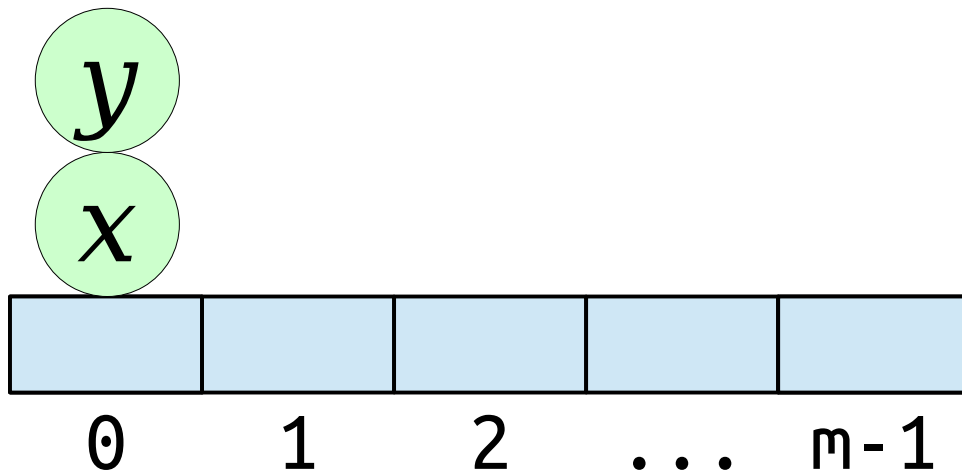$$= \sum_{i=0}^{m-1} \frac{1}{m^2}$$

$y$

$x$

0   1   2   ...   m-1

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

$$= \sum_{i=0}^{m-1} \frac{1}{m^2}$$

$$= \frac{1}{m}$$

$y$

$x$

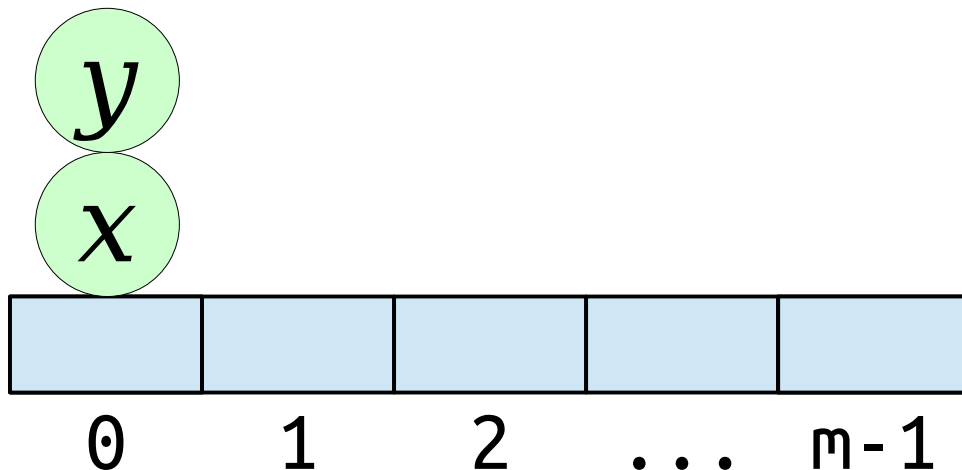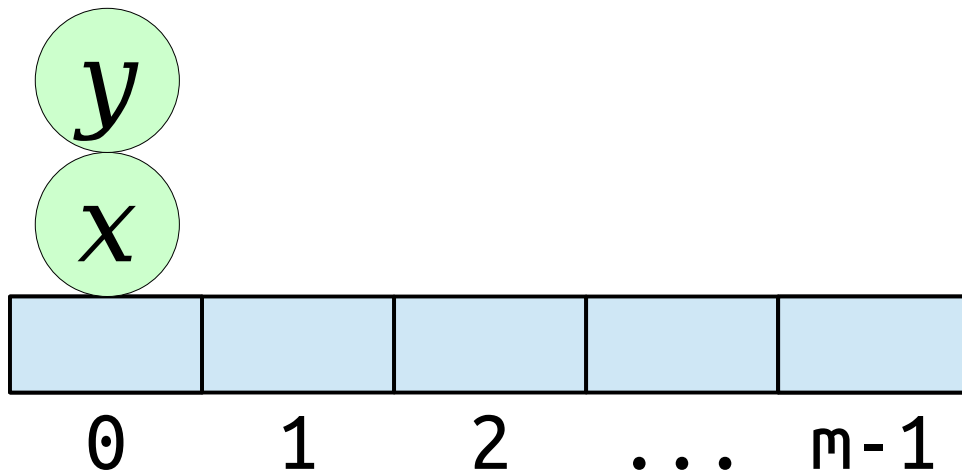| | | | | |
|---|---|---|---|---|
| | | | | |

0    1    2    ...   m-1

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

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

**Intuition:**
2-independence means any pair of elements is unlikely to collide.

$$\Pr[h(x) = h(y)]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i \wedge h(y) = i]$$

$$= \sum_{i=0}^{m-1} \Pr[h(x) = i] \cdot \Pr[h(y) = i]$$

$$= \sum_{i=0}^{m-1} \frac{1}{m^2}$$

$$= \frac{1}{m}$$

This is the same as if $h$ were a truly random function.

$y$

$x$

0   1   2   ...   m-1

For more on hashing outside of Theoryland, check out **_this Stack Exchange post_**.

# Approximating Quantities

# What makes for a good "approximate" solution?

Let **A** be the true answer. Let **Â** be a random variable denoting our estimate.

This would not make for a good estimate. However, we have E[Â] = A.

**Observation 1:** Being correct in expectation isn't sufficient.

**Distribution of our estimate Â.**

**A (True answer)**

# What does it mean for an approximation to be "good"?

Let **A** be the true answer. Let **Â** be a random variable denoting our estimate.

It's unlikely that we'll get the right answer, but we're probably going to be close.

***Observation 2:*** The difference $|\hat{A} - A|$ between our estimate and the truth should ideally be small.

Distribution of our estimate Â.

**A**
*(True answer)*

# What does it mean for an approximation to be "good"?

Let **A** be the true answer. Let **Â** be a random variable denoting our estimate.

This estimate skews low, but it's very close to the true value.

***Observation 3:*** An estimate doesn't have to be unbiased to be useful.



*Distribution of our estimate Â.*

**A**
*(True answer)*

# What does it mean for an approximation to be "good"?

Let $A$ be the true answer. Let $\hat{A}$ be a random variable denoting our estimate.

**Memory used: 16MB**



$A$
*(True answer)*

# What does it mean for an approximation to be "good"?

Let **A** be the true answer. Let **Â** be a random variable denoting our estimate.

**Memory used: 32MB**



**A**
*(True answer)*

# What does it mean for an approximation to be "good"?

Let **A** be the true answer. Let $\hat{A}$ be a random variable denoting our estimate.

**Memory used: 64MB**



**A**
*(True answer)*

# What does it mean for an approximation to be "good"?

Let **A** be the true answer. Let **Â** be a random variable denoting our estimate.

The more resources we allocate, the better our estimate should be.

***Observation 4:*** A good approximation should be tunable.

**Memory used: 256MB**

**A**
*(True answer)*

# What does it mean for an approximation to be "good"?

We have two user-provided values

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

where $\varepsilon$ represents **_accuracy_** and $\delta$ represents **_confidence_**.

**_Goal:_** Make an estimator $\hat{A}$ for some quantity $A$ where

With probability at least $1 - \delta$,

$$|\hat{A} - A| \leq \varepsilon \cdot size(input)$$

for some measure of the size of the input.

What does it mean for an approximation to be "good"?

We have two user-provided values

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

where $\varepsilon$ represents **accuracy** and $\delta$ represents **confidence**.

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

With probability at least $1 - \delta$, $\quad\longleftarrow$ — ***Probably***

$$|\hat{A} - A| \leq \varepsilon \cdot size(input)$$

for some measure of the size of the input.

What does it mean for an approximation to be "good"?

We have two user-provided values

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

where $\varepsilon$ represents **accuracy** and $\delta$ represents **confidence**.

---

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

With probability at least $1 - \delta$,  ⟵ ────── **Probably**

$$|\hat{A} - A| \leq \varepsilon \cdot size(input)$$  ⟵ ────── **Approximately Correct**

for some measure of the size of the input.

---

# What does it mean for an approximation to be "good"?

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

With probability at least 1 – δ,  ⎫⎬⎭ ← ***Probably***

$$|A - \hat{A}| \leq \varepsilon \cdot size(input)$$  ⎫⎬⎭ ← ***Approximately Correct***

for some measure of the size of the input.

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta$, ⟶ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟶ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ small

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta,$ ⟵ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot \mathit{size}(\mathit{input})$ ⟵ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ small

What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta$, ⟵ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ small



**True answer**

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta$, ⟵ ── *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟵ ── *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ small

*True answer*

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta$, ⟸ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟸ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ small

*True answer*

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta$, ⟵ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ medium

*True answer*



# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta$, ⟵ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ large

**True answer**

What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta,$ ⟶ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟶ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{2}$
$\varepsilon$ small

*True answer*

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta,$ ⟵ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of the size of the input.

$\delta = \frac{1}{4}$
$\varepsilon$ small

*True answer*

# What does it mean for an approximation to be "good"?

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

With probability at least $1 - \delta,$ ⟵ *Probably*

$|A - \hat{A}| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of the size of the input.

$\delta = {}^{1}/_{16}$
$\varepsilon$ small

*True answer*

# What does it mean for an approximation to be "good"?

# Frequency Estimation

# Frequency Estimators

- A ***frequency estimator*** is a data structure supporting the following operations:

  - ***increment***$(x)$, which increments the number of times that $x$ has been seen, and

  - ***estimate***$(x)$, which returns an estimate of the frequency of $x$.

- Using BSTs, we can solve this in space $\Theta(n)$ with worst-case $O(\log n)$ costs on the operations.

- Using hash tables, we can solve this in space $\Theta(n)$ with expected $O(1)$ costs on the operations.

# Frequency Estimators

- Frequency estimation has many applications:

  - Search engines: Finding frequent search queries.

  - Network routing: Finding common source and destination addresses.

- In these applications, $\Theta(n)$ memory can be impractical.

- **Goal:** Get *approximate* answers to these queries in sublinear space.

# The Count-Min Sketch

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. Use a ***sum of indicator variables*** and ***linearity of expectation*** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a ***concentration inequality*** to show that, with decent probability, the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

# How to Build an Estimator

1. **Design a simple data structure that, intuitively, gives you a good estimate.**

2. Use a **sum of indicator variables** and **linearity of expectation** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a **concentration inequality** to show that, with decent probability, the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

# Revisiting the Exact Solution

- In the exact solution to the frequency estimation problem, we maintained a single counter for each distinct element. This is too space-inefficient.

- **Idea:** Store a fixed number of counters and assign a counter to each $x_i \in \mathcal{U}$. Multiple $x_i$'s might be assigned to the same counter.

- To **increment**$(x)$, increment the counter for $x$.

- To **estimate**$(x)$, read the value of the counter for $x$.

# Revisiting the Exact Solution

- In the exact solution to the frequency estimation problem, we maintained a single counter for each distinct element. This is too space-inefficient.

- *Idea:* Store a fixed number of counters and assign a counter to each $x_i \in \mathcal{U}$. Multiple $x_i$'s might be assigned to the same counter.

- To *increment*$(x)$, increment the counter for $x$.

- To *estimate*$(x)$, read the value of the counter for $x$.

# Revisiting the Exact Solution

- In the exact solution to the frequency estimation problem, we maintained a single counter for each distinct element. This is too space-inefficient.

- **_Idea:_** Store a fixed number of counters and assign a counter to each $x_i \in \mathcal{U}$. Multiple $x_i$'s might be assigned to the same counter.

- To **_increment_**$(x)$, increment the counter for $x$.

- To **_estimate_**$(x)$, read the value of the counter for $x$.



| 12 | 6 | 5 | 7 |

# Our Initial Structure

- We can model "assigning each $x_i$ to a counter" by using hash functions.

- Choose, from a family of 2-independent hash functions $\mathscr{H}$, a uniformly-random hash function $h : \mathscr{U} \to [w]$.

- Create an array **count** of $w$ counters, each initially zero.

  - We'll choose $w$ later on.

- To *increment*($x$), increment **count**$[h(x)]$.

- To *estimate*($x$), return **count**$[h(x)]$.

# Analyzing our Structure

For each $x_i \in \mathcal{U}$, let $\boldsymbol{a}_i$ denote the number of times we've seen $x_i$.

Similarly, let $\hat{\boldsymbol{a}}_i$ denote our estimated value of the frequency of $x_i$.

***Goal:*** Bound the probability that the error $(\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i)$ is too high.

**Idea:** Think of our element frequencies $a_1$, $a_2$, $a_3$, ... as a vector

$$a = [a_1, a_2, a_3, ... ].$$

The total number of objects is the sum of the vector entries.

This is called the **$L_1$ norm** of $a$, and is denoted $\|a\|_1$:

$$\|a\|_1 = \sum_i |a_i|$$

There are $\|a\|_1$ total elements distributed across $w$ buckets. We're using a 2-independent hash family.

**Reasonable guess:** each bin has $\|a\|_1 / w$ elements in it, so
$$\hat{a}_i - a_i \leq \|a\|_1 / w$$



9       5       4

**Number of buckets: $w$**

**Question:** Intuitively, what should we expect our approximation error to be?

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. Use a ***sum of indicator variables*** and ***linearity of expectation*** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a ***concentration inequality*** to show that, with decent probability, the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. **Use a *sum of indicator variables* and *linearity of expectation* to prove that, on expectation, the data structure is pretty close to correct.**

3. Use a ***concentration inequality*** to show that, with decent probability, the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

# Analyzing this Structure

- Let's look at $\hat{a}_i = \mathbf{\color{purple}{count}}[h(x_i)]$ for some choice of $x_i$.

- For each element $x_j$:

    - If $h(x_i) = h(x_j)$, then $x_j$ contributes $a_j$ to $\mathbf{\color{purple}{count}}[h(x_i)]$.

    - If $h(x_i) \neq h(x_j)$, then $x_j$ contributes $0$ to $\mathbf{\color{purple}{count}}[h(x_i)]$.

# Analyzing this Structure

- Let's look at $\hat{a}_i =$ **count**$[h(x_i)]$ for some choice of $x_i$.

- For each element $x_j$:

  - If $h(x_i) = h(x_j)$, then $x_j$ contributes $a_j$ to **count**$[h(x_i)]$.

  - If $h(x_i) \neq h(x_j)$, then $x_j$ contributes $0$ to **count**$[h(x_i)]$.

- To pin this down precisely, let's define a set of random variables $X_1, X_2, \ldots,$ as follows:

$$X_j = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$$

# Analyzing this Structure

- Let's look at $\hat{a}_i = \mathbf{count}[h(x_i)]$ for some choice of $x_i$.

- For each element $x_j$:

  - If $h(x_i) = h(x_j)$, then $x_j$ contributes $a_j$ to $\mathbf{count}[h(x_i)]$.

  - If $h(x_i) \neq h(x_j)$, then $x_j$ contributes 0 to $\mathbf{count}[h(x_i)]$.

- To pin this down precisely, let's define a set of random variables $X_1, X_2, \ldots,$ as follows:

$$X_j = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$$

Each of these variables is called an ***indicator random variable***, since it "indicates" whether some event occurs.

# Analyzing this Structure

- Let's look at $\hat{a}_i = \textbf{count}[h(x_i)]$ for some choice of $x_i$.
- For each element $x_j$:
  - If $h(x_i) = h(x_j)$, then $x_j$ contributes $a_j$ to $\textbf{count}[h(x_i)]$.
  - If $h(x_i) \neq h(x_j)$, then $x_j$ contributes $0$ to $\textbf{count}[h(x_i)]$.
- To pin this down precisely, let's define a set of random variables $X_1, X_2, \ldots$, as follows:

$$X_j = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$$

- The value of $\hat{a}_i - a_i$ is then given by

$$\hat{a}_i - a_i = \sum_{j \neq i} a_j X_j$$

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] \;=\; \mathrm{E}\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] = \mathrm{E}\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \mathrm{E}\left[\boldsymbol{a}_j X_j\right]$$

$$\mathrm{E}[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] \;=\; \mathrm{E}\Big[\sum_{j \neq i} \boldsymbol{a}_j X_j\Big]$$

$$=\; \sum_{j \neq i} \mathrm{E}[\boldsymbol{a}_j X_j]$$

This follows from ***linearity of expectation***. We'll use this property extensively over the next few days.

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] = \mathrm{E}\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \mathrm{E}\left[\boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}\left[X_j\right]$$

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] \;=\; \mathrm{E}\!\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$=\; \sum_{j \neq i} \mathrm{E}\!\left[\boldsymbol{a}_j X_j\right]$$

$$=\; \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}\!\left[X_j\right]$$

The values of $\boldsymbol{a}_j$ are not random. ***The randomness comes from our choice of hash function.***

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] \;=\; \mathrm{E}\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$=\; \sum_{j \neq i} \mathrm{E}\left[\boldsymbol{a}_j X_j\right]$$

$$=\; \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}\left[X_j\right]$$

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = E\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} E\left[\boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j E\left[X_j\right]$$

---

$$E\left[X_j\right] =$$

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = E[\sum_{j \neq i} \boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} E[\boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j E[X_j]$$

---

$$E[X_j] =$$

$$X_j = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$$

$$\mathrm{E}[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] \;=\; \mathrm{E}\Big[\sum_{j \neq i} \boldsymbol{a}_j X_j\Big]$$

$$=\; \sum_{j \neq i} \mathrm{E}[\boldsymbol{a}_j X_j]$$

$$=\; \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}[X_j]$$

---

$$\mathrm{E}[X_j] \;=\; 1 \cdot \mathrm{Pr}[h(x_i)=h(x_j)] \;+\; 0 \cdot \mathrm{Pr}[h(x_i) \neq h(x_j)]$$

$$X_j = \begin{cases} 1 & \text{if } h(x_i)=h(x_j) \\ 0 & \text{otherwise} \end{cases}$$

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = E\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} E\left[\boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j E\left[X_j\right]$$

---

$$E\left[X_j\right] = 1 \cdot \Pr\left[h(x_i) = h(x_j)\right] + 0 \cdot \Pr\left[h(x_i) \neq h(x_j)\right]$$

$$\mathrm{E}[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] \;=\; \mathrm{E}\Big[\sum_{j \neq i} \boldsymbol{a}_j X_j\Big]$$

$$=\; \sum_{j \neq i} \mathrm{E}[\boldsymbol{a}_j X_j]$$

$$=\; \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}[X_j]$$

---

$$\mathrm{E}[X_j] \;=\; 1 \cdot \Pr[h(x_i) = h(x_j)] \;+\; 0 \cdot \Pr[h(x_i) \neq h(x_j)]$$

$$=\; \Pr[h(x_i) = h(x_j)]$$

$$\mathrm{E}[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = \mathrm{E}\Big[\sum_{j \neq i} \boldsymbol{a}_j X_j\Big]$$

$$= \sum_{j \neq i} \mathrm{E}[\boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}[X_j]$$

---

$$\mathrm{E}[X_j] = 1 \cdot \mathrm{Pr}[h(x_i) = h(x_j)] + 0 \cdot \mathrm{Pr}[h(x_i) \neq h(x_j)]$$

$$= \mathrm{Pr}[h(x_i) = h(x_j)]$$

If $X$ is an indicator variable for some event $\mathcal{E}$, then **E[X] = Pr[$\mathcal{E}$]**. This is really useful when using linearity of expectation!

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = E\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} E[\boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j E[X_j]$$

---

$$E[X_j] = 1 \cdot \Pr[h(x_i) = h(x_j)] + 0 \cdot \Pr[h(x_i) \neq h(x_j)]$$

$$= \Pr[h(x_i) = h(x_j)]$$

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] \;=\; E\!\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$=\; \sum_{j \neq i} E[\boldsymbol{a}_j X_j]$$

$$=\; \sum_{j \neq i} \boldsymbol{a}_j E[X_j]$$

---

$$E[X_j] \;=\; 1 \cdot \Pr[h(x_i) = h(x_j)] \;+\; 0 \cdot \Pr[h(x_i) \neq h(x_j)]$$

$$=\; \Pr[h(x_i) = h(x_j)]$$

> Hey, we saw this earlier!

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = E[\sum_{j \neq i} \boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} E[\boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j E[X_j]$$

---

$$E[X_j] = 1 \cdot Pr[h(x_i) = h(x_j)] + 0 \cdot Pr[h(x_i) \neq h(x_j)]$$

$$= Pr[h(x_i) = h(x_j)]$$

$$= \frac{1}{w}$$

Hey, we saw this earlier!

$$E[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = E\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} E[\boldsymbol{a}_j X_j]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j E[X_j]$$

$$= \sum_{j \neq i} \frac{\boldsymbol{a}_j}{w}$$

---

$$E[X_j] = 1 \cdot \Pr[h(x_i) = h(x_j)] + 0 \cdot \Pr[h(x_i) \neq h(x_j)]$$

$$= \Pr[h(x_i) = h(x_j)]$$

$$= \frac{1}{w}$$

Hey, we saw this earlier!

$$\mathrm{E}[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] \;=\; \mathrm{E}\Big[\sum_{j \neq i} \boldsymbol{a}_j X_j\Big]$$

$$=\; \sum_{j \neq i} \mathrm{E}[\boldsymbol{a}_j X_j]$$

$$=\; \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}[X_j]$$

$$=\; \sum_{j \neq i} \frac{\boldsymbol{a}_j}{w}$$

---

$$\mathrm{E}[X_j] \;=\; 1 \cdot \Pr[h(x_i)=h(x_j)] \;+\; 0 \cdot \Pr[h(x_i) \neq h(x_j)]$$
$$=\; \Pr[h(x_i)=h(x_j)]$$
$$=\; \frac{1}{w}$$

$$\mathrm{E}[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i] = \mathrm{E}\left[\sum_{j \neq i} \boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \mathrm{E}\left[\boldsymbol{a}_j X_j\right]$$

$$= \sum_{j \neq i} \boldsymbol{a}_j \mathrm{E}\left[X_j\right]$$

$$= \sum_{j \neq i} \frac{\boldsymbol{a}_j}{w}$$

$$\leq \frac{\|\boldsymbol{a}\|_1}{w}$$

---

$$\mathrm{E}[X_j] = 1 \cdot \Pr[h(x_i) = h(x_j)] + 0 \cdot \Pr[h(x_i) \neq h(x_j)]$$

$$= \Pr[h(x_i) = h(x_j)]$$

$$= \frac{1}{w}$$

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. Use a ***sum of indicator variables*** and ***linearity of expectation*** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a ***concentration inequality*** to show that, with decent probability, the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. Use a ***sum of indicator variables*** and ***linearity of expectation*** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a ***concentration inequality*** to show that, with decent probability, the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

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

With probability at least $1 - \delta$, ⟵ ─── *Probably*

$|\hat{a} - a| \leq \varepsilon \cdot size(input)$ ⟵ ─── *Approximately Correct*

for some measure of the size of the input.

$\varepsilon \|a\|_1$

$a_i$

$$\mathrm{E}\left[\hat{a}_i - a_i\right] \leq \frac{\|a\|_1}{w}$$

How do we tune $w$ so we're likely to fall in this range?

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \ > \ \varepsilon\left\|\boldsymbol{a}\right\|_1\right]$$

$$\Pr\left[\boxed{\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i} \; > \; \varepsilon\left\|\boldsymbol{a}\right\|_1\right]$$

$$\Pr\left[\; \boxed{\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i} \; > \; \varepsilon \left\|\boldsymbol{a}\right\|_1\right]$$

We don't know the exact distribution of this random variable.

However, we have a ***one-sided error***: our estimate can never be lower than the true value. This means that $\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \geq 0$.

***Markov's inequality*** says that if $X$ is a nonnegative random variable, then

$$\Pr[X \geq c] \; \leq \; \frac{\mathrm{E}[X]}{c}.$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \; > \; \varepsilon \, \|\boldsymbol{a}\|_1\right]$$

$$\leq \quad \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon \, \|\boldsymbol{a}\|_1}$$

We don't know the exact distribution of this random variable.

However, we have a **_one-sided error_**: our estimate can never be lower than the true value. This means that $\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \geq 0$.

**_Markov's inequality_** says that if $X$ is a nonnegative random variable, then

$$\Pr[X \geq c] \; \leq \; \frac{\mathrm{E}[X]}{c}.$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \; > \; \varepsilon \left\|\boldsymbol{a}\right\|_1\right]$$

$$\leq \quad \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon \left\|\boldsymbol{a}\right\|_1}$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \; > \; \varepsilon\,\|\boldsymbol{a}\|_1\right]$$

$$\leq \quad \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon\,\|\boldsymbol{a}\|_1}$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \; > \; \varepsilon \left\|\boldsymbol{a}\right\|_1\right]$$

$$\leq \quad \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon \left\|\boldsymbol{a}\right\|_1}$$

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] \; \leq \; \frac{\left\|\boldsymbol{a}\right\|_1}{w}$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \; > \; \varepsilon \, \|\boldsymbol{a}\|_1\right]$$

$$\leq \quad \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon \, \|\boldsymbol{a}\|_1}$$

$$\leq \quad \frac{\|\boldsymbol{a}\|_1}{w} \cdot \frac{1}{\varepsilon \, \|\boldsymbol{a}\|_1}$$

$$\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right] \; \leq \; \frac{\|\boldsymbol{a}\|_1}{w}$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1\right]$$

$$\leq \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon \|\boldsymbol{a}\|_1}$$

$$\leq \frac{\|\boldsymbol{a}\|_1}{w} \cdot \frac{1}{\varepsilon \|\boldsymbol{a}\|_1}$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \ > \ \varepsilon \left\|\boldsymbol{a}\right\|_1\right]$$

$$\leq \quad \frac{\mathrm{E}\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i\right]}{\varepsilon \left\|\boldsymbol{a}\right\|_1}$$

$$\leq \quad \frac{\left\|\boldsymbol{a}\right\|_1}{w} \cdot \frac{1}{\varepsilon \left\|\boldsymbol{a}\right\|_1}$$

$$= \quad \frac{1}{\varepsilon \, w}$$

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

With probability at least $1 - \delta$, ⟵ *Probably*

$|\hat{\boldsymbol{a}} - \boldsymbol{a}| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of input size.

$$\Pr[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1] \leq \frac{1}{\varepsilon w}$$

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

With probability at least $1 - \delta,$  ⟵ *Probably*

$|\hat{a} - a| \leq \varepsilon \cdot size(input)$  ⟵ *Approximately Correct*

for some measure of input size.

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \ \leq \ \frac{1}{\varepsilon\, w}$$

**Initial Idea:**

Pick $w = \varepsilon^{-1} \cdot \delta^{-1}$. Then

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \ \leq \ \delta$$

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

With probability at least $1 - \delta$, ⟵ *Probably*

$|\hat{a} - a| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of input size.

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \leq \frac{1}{\varepsilon w}$$

Suppose we're counting 1,000 distinct items.

**Initial Idea:**
Pick $w = \varepsilon^{-1} \cdot \delta^{-1}$. Then

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \leq \delta$$

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

With probability at least $1 - \delta,$ — *Probably*

$|\hat{a} - a| \leq \varepsilon \cdot size(input)$ — *Approximately Correct*

for some measure of input size.

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \leq \frac{1}{\varepsilon w}$$

Suppose we're counting 1,000 distinct items.

If we want our estimate to be within $\varepsilon\|a\|_1$ of the true value with 99.9% probability, how much memory do we need?

**Initial Idea:**
Pick $w = \varepsilon^{-1} \cdot \delta^{-1}$. Then

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \leq \delta$$

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

With probability at least $1 - \delta$, — *Probably*

$|\hat{\boldsymbol{a}} - \boldsymbol{a}| \leq \varepsilon \cdot size(input)$ — *Approximately Correct*

for some measure of input size.

$$\Pr[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1] \leq \frac{1}{\varepsilon w}$$

**Initial Idea:**
Pick $w = \varepsilon^{-1} \cdot \delta^{-1}$. Then

$$\Pr[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1] \leq \delta$$

Suppose we're counting 1,000 distinct items.

If we want our estimate to be within $\varepsilon \|\boldsymbol{a}\|_1$ of the true value with 99.9% probability, how much memory do we need?

***Answer:*** **1,000** $\cdot$ **$\varepsilon^{-1}$**.

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

With probability at least $1 - \delta$, ⟶ *Probably*

$|\hat{a} - a| \le \varepsilon \cdot size(input)$ ⟶ *Approximately Correct*

for some measure of input size.

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \le \frac{1}{\varepsilon w}$$

**Initial Idea:**
Pick $w = \varepsilon^{-1} \cdot \delta^{-1}$. Then

$$\Pr[\hat{a}_i - a_i > \varepsilon \|a\|_1] \le \delta$$

Suppose we're counting 1,000 distinct items.

If we want our estimate to be within $\varepsilon \|a\|_1$ of the true value with 99.9% probability, how much memory do we need?

***Answer:*** **1,000** $\cdot$ **$\varepsilon^{-1}$**.

***Can we do better?***

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. Use a ***sum of indicator variables*** and ***linearity of expectation*** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a ***concentration inequality*** to show that the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

# How to Build an Estimator

1. Design a simple data structure that, intuitively, gives you a good estimate.

2. Use a **sum of indicator variables** and **linearity of expectation** to prove that, on expectation, the data structure is pretty close to correct.

3. Use a **concentration inequality** to show that the data structure's output is close to its expectation.

4. Run multiple copies of the data structure in parallel to amplify the success probability.

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

With probability at least $1 - \delta$, — *Probably*

$|\hat{\boldsymbol{a}} - \boldsymbol{a}| \leq \varepsilon \cdot size(input)$ — *Approximately Correct*

for some measure of input size.

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1\right] \leq \frac{1}{\varepsilon w}$$

We could choose $w = k \cdot \varepsilon^{-1}$ for any constant $k$ to get a failure probability of at most $k^{-1}$. The choice of $e$ is (mostly) arbitrary.

**Revised Idea:** Pick $w = e \cdot \varepsilon^{-1}$. Then

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1\right] < e^{-1}$$

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

With probability at least $1 - \delta$, $\}$ ⟵ *Probably*

$|\hat{a} - a| \le \varepsilon \cdot size(input)$ $\}$ ⟵ *Approximately Correct*

for some measure of input size.

$$\Pr\left[\hat{a}_i - a_i > \varepsilon \|a\|_1\right] \le \frac{1}{\varepsilon w}$$

This simple data structure, by itself, is likely to be wrong.

What happens if we run a bunch of copies of this approach in parallel?

**Revised Idea:** Pick $w = e \cdot \varepsilon^{-1}$. Then

$$\Pr\left[\hat{a}_i - a_i > \varepsilon \|a\|_1\right] < e^{-1}$$

# Running in Parallel

- Let's run **d** copies of our data structure in parallel with one another.

- Each row has its hash function sampled uniformly at random from our hash family.

- Each time we *increment* an item, we perform the corresponding *increment* operation on each row.

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

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 31 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 39 | 88 | ... | 75 |
| ... | | | | ... | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 59 |

$d = ??$

# Running in Parallel

- Let's run **d** copies of our data structure in parallel with one another.

- Each row has its hash function sampled uniformly at random from our hash family.

- Each time we *increment* an item, we perform the corresponding *increment* operation on each row.

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

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 31 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 39 | 88 | ... | 75 |
| ... | | | | ... | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 59 |

$d = ??$

# Running in Parallel

- Let's run **d** copies of our data structure in parallel with one another.

- Each row has its hash function sampled uniformly at random from our hash family.

- Each time we *increment* an item, we perform the corresponding *increment* operation on each row.

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

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 32 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 29 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 40 | 88 | ... | 75 |
| ... | ... | | | | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 60 |

$d = ??$

# Running in Parallel

- Let's run **_d_** copies of our data structure in parallel with one another.

- Each row has its hash function sampled uniformly at random from our hash family.

- Each time we **_increment_** an item, we perform the corresponding **_increment_** operation on each row.

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

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 32 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 29 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 40 | 88 | ... | 75 |
| ... | | | | ... | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 60 |

$d = ??$

# Running in Parallel

- Imagine we 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?

Formulate a hypothesis, but ***don't post anything in chat just yet***.

| *Estimator 1:* | *Estimator 2:* | *Estimator 3:* | *Estimator 4:* | *Estimator 5:* |
|---|---|---|---|---|
| 137 | 271 | 166 | 103 | 261 |

# Running in Parallel

- Imagine we 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?

> Now, **private chat me your best guess**. Not sure? Just answer "??".

| Estimator 1: | Estimator 2: | Estimator 3: | Estimator 4: | Estimator 5: |
|:---:|:---:|:---:|:---:|:---:|
| 137 | 271 | 166 | 103 | 261 |

# Running in Parallel

- Imagine we 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?

| Estimator 1: | Estimator 2: | Estimator 3: | Estimator 4: | Estimator 5: |
|:---:|:---:|:---:|:---:|:---:|
| 137 | 271 | 166 | 103 | 261 |

# Running in Parallel

- Imagine we 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?

| *Estimator 1:* | *Estimator 2:* | *Estimator 3:* | *Estimator 4:* | *Estimator 5:* |
|:---:|:---:|:---:|:---:|:---:|
| 137 | 271 | 166 | **103** | 261 |

# Running in Parallel

- Imagine we 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 aggr into a single estimate?

**Intuition:** The smallest estimate returned has the least "noise," and that's the best guess for the frequency.

| *Estimator 1:* | *Estimator 2:* | *Estimator 3:* | *Estimator 4:* | *Estimator 5:* |
|:---:|:---:|:---:|:---:|:---:|
| 137 | 271 | 166 | 103 | 261 |

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is min $\{\hat{\boldsymbol{a}}_{ij}\}$

$$\Pr\left[\min\left\{\,\hat{\boldsymbol{a}}_{ij}\,\right\}-\boldsymbol{a}_i\ >\ \varepsilon\,\|\boldsymbol{a}\|_1\right]$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\{\hat{\boldsymbol{a}}_{ij}\}$

$$\Pr \left[ \min \left\{ \, \hat{\boldsymbol{a}}_{ij} \, \right\} - \boldsymbol{a}_i \; > \; \varepsilon \left\| \boldsymbol{a} \right\|_1 \right]$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min \left\{ \hat{\boldsymbol{a}}_{ij} \right\}$

$$\Pr\left[\min\left\{\hat{\boldsymbol{a}}_{ij}\right\} - \boldsymbol{a}_i \ > \ \varepsilon\left\|\boldsymbol{a}\right\|_1\right]$$

The only way the minimum estimate is inaccurate is if *every* estimate is inaccurate.

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is min $\{\hat{\boldsymbol{a}}_{ij}\}$

$$\Pr\left[\;\boxed{\min\{\,\hat{a}_{ij}\,\} - a_i \;>\; \varepsilon\,\|a\|_1}\;\right]$$

$$=\;\Pr\left[\;\bigwedge_{j=1}^{d}\left(\hat{a}_{ij} - a_i \;>\; \varepsilon\,\|a\|_1\right)\right]$$

The only way the minimum estimate is inaccurate is if *every* estimate is inaccurate.

Let $\hat{a}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\{\hat{a}_{ij}\}$

$$\Pr\left[\min\left\{\hat{\boldsymbol{a}}_{ij}\right\} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$= \Pr\left[\bigwedge_{j=1}^{d}\left(\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right)\right]$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\left\{\hat{\boldsymbol{a}}_{ij}\right\}$

$$\Pr\left[\min\{\hat{\boldsymbol{a}}_{ij}\} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$= \Pr\left[\bigwedge_{j=1}^{d}\left(\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right)\right]$$

Each copy of the data structure is independent of the others.

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\{\hat{\boldsymbol{a}}_{ij}\}$

$$\Pr\left[\min\left\{\hat{\boldsymbol{a}}_{ij}\right\} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$= \Pr\left[\bigwedge_{j=1}^{d}\left(\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right)\right]$$

$$= \prod_{j=1}^{d}\Pr\left[\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

Each copy of the data structure is independent of the others.

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\left\{\hat{\boldsymbol{a}}_{ij}\right\}$

$$\Pr\left[\min\left\{\hat{\boldsymbol{a}}_{ij}\right\} - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1\right]$$

$$= \Pr\left[\bigwedge_{j=1}^{d}\left(\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1\right)\right]$$

$$= \prod_{j=1}^{d} \Pr\left[\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon \|\boldsymbol{a}\|_1\right]$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\left\{\hat{\boldsymbol{a}}_{ij}\right\}$

$$\Pr \left[ \min \left\{ \hat{\boldsymbol{a}}_{ij} \right\} - \boldsymbol{a}_i \; > \; \varepsilon \left\| \boldsymbol{a} \right\|_1 \right]$$

$$= \; \Pr \left[ \bigwedge_{j=1}^{d} \left( \hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i \; > \; \varepsilon \left\| \boldsymbol{a} \right\|_1 \right) \right]$$

$$= \; \prod_{j=1}^{d} \Pr \left[ \hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i \; > \; \varepsilon \left\| \boldsymbol{a} \right\|_1 \right]$$

$$\Pr \left[ \hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \; \geq \; \varepsilon \left\| \boldsymbol{a} \right\|_1 \right] \; \leq \; e^{-1}$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min \left\{ \hat{\boldsymbol{a}}_{ij} \right\}$

$$\Pr\left[\min\left\{\hat{\boldsymbol{a}}_{ij}\right\} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$= \Pr\left[\bigwedge_{j=1}^{d}\left(\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right)\right]$$

$$= \prod_{j=1}^{d}\Pr\left[\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$\leq \prod_{j=1}^{d} e^{-1}$$

$$\Pr\left[\hat{\boldsymbol{a}}_i - \boldsymbol{a}_i \geq \varepsilon\|\boldsymbol{a}\|_1\right] \leq e^{-1}$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\{\hat{\boldsymbol{a}}_{ij}\}$

$$\Pr\left[\min\left\{\hat{\boldsymbol{a}}_{ij}\right\} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$= \Pr\left[\bigwedge_{j=1}^{d}\left(\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right)\right]$$

$$= \prod_{j=1}^{d}\Pr\left[\hat{\boldsymbol{a}}_{ij} - \boldsymbol{a}_i > \varepsilon\|\boldsymbol{a}\|_1\right]$$

$$\leq \prod_{j=1}^{d}e^{-1}$$

$$= e^{-d}$$

Let $\hat{\boldsymbol{a}}_{ij}$ be the estimate from the $j$th copy of the data structure.

Our final estimate is $\min\left\{\hat{\boldsymbol{a}}_{ij}\right\}$

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

With probability at least $1 - \delta$, ⟵ ⟵ *Probably*

$|\hat{a} - a| \leq \varepsilon \cdot size(input)$ ⟵ *Approximately Correct*

for some measure of input size.

$$\Pr\left[\min\{\hat{a}_{ij}\} - a_i > \varepsilon\|a\|_1\right] \leq e^{-d}$$

**Idea:** Choose $d = -\ln \delta$. (Equivalently: $d = \ln \delta^{-1}$.) Then

$$\Pr\left[\min\{\hat{a}_{ij}\} - a_i > \varepsilon\|a\|_1\right] \leq \delta$$

# The Count-Min Sketch

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

| | | | | | | |
|---|---|---|---|---|---|---|
| $h_1$ | 31 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 39 | 88 | ... | 75 |
| ... | ... |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 59 |

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

Sampled uniformly and independently from a 2-independent family of hash functions

# The Count-Min Sketch

| $h_1$ | 31 | 41 | 59 | 26 | 53 | ... | 58 |
|-------|----|----|----|----|----|-----|----|
| $h_2$ | 27 | 18 | 28 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 39 | 88 | ... | 75 |
| ... | | | ... | | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 59 |

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

# The Count-Min Sketch



| $h_1$ | 31 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 18 | 28 | ... | 45 |
| $h_3$ | 16 | 18 | 3 | 39 | 88 | ... | 75 |
| *...* | | | | ... | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 59 |

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

# The Count-Min Sketch

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 32 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 19 | 28 | ... | 45 |
| $h_3$ | 16 | 19 | 3 | 39 | 88 | ... | 75 |
| ... | | | | ... | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 60 |

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

# The Count-Min Sketch

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 32 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 19 | 28 | ... | 45 |
| $h_3$ | 16 | 19 | 3 | 39 | 88 | ... | 75 |
| ... | ... | | | | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 60 |

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

# The Count-Min Sketch

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $h_1$ | 32 | 41 | 59 | 26 | 53 | ... | 58 |
| $h_2$ | 27 | 18 | 28 | 19 | 28 | ... | 45 |
| $h_3$ | 16 | 19 | 3 | 39 | 88 | ... | 75 |
| ... | ... | | | | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 60 |

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

```
estimate(x):
    result = ∞
    for i = 1 … d:
        result = min(result, count[i][hᵢ(x)])
    return result
```

# The Count-Min Sketch

| $h_1$ | 32 | 41 | 59 | 26 | 53 | ... | 58 |
|---|---|---|---|---|---|---|---|
| $h_2$ | 27 | 18 | 28 | 19 | 28 | ... | 45 |
| $h_3$ | 16 | 19 | 3 | 39 | 88 | ... | 75 |
| ... | | | | ... | | | |
| $h_d$ | 69 | 31 | 47 | 18 | 5 | ... | 60 |

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

```
estimate(x):
    result = ∞
    for i = 1 … d:
        result = min(result, count[i][hᵢ(x)])
    return result
```

# The Count-Min Sketch

- Update and query times are $\Theta(d)$, which is $\Theta(\log \delta^{-1})$.

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

  - This is a *major* improvement over our earlier approach that used $\Theta(\varepsilon^{-1} \cdot \delta^{-1})$ counters.

  - This can be *significantly* better than just storing a raw frequency count!

- Provides an estimate to within $\varepsilon \|a\|_1$ with probability at least $1 - \delta$.

# Major Ideas From Today

- ***2-independent hash families*** are useful when we want to keep collisions low.

- A "good" approximation of some quantity should have tunable ***confidence*** and ***accuracy*** parameters.

- ***Sums of indicator variables*** are useful for deriving expected values of estimators.

- ***Concentration inequalities*** like ***Markov's inequality*** are useful for showing estimators don't stay too much from their expected values.

- Good estimators can be built from multiple parallel copies of weaker estimators.

# Next Time

- ***Count Sketches***

  - An alternative frequency estimator with different time/space bounds.

- ***Cardinality Estimation***

  - Estimating how many different items you've seen in a data stream.