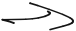# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.

Lecture 3

## Logistics

- Problem Set 1 has been posted on the course website and is due Friday 9/22 at 11:59pm.

- On the quiz feedback question, several people mentioned concerns about linear algebra background. The good news is you have some time – we will do essentially no linear algebra before the midterm. See resources on Lecture 15 for review material to get started on.

- For probability/problem solving practice beyond the quizzes/pset I highly recommend looking at the exercises in *Foundations of Data Science* and *Probability and Computing*. Feel free to ask for solutions to these on Piazza.

- It is common to not catch everything in lecture. I strongly encourage going back to the slides to review/check your understanding after class. Also come to office hours for more in-depth discussion/examples.

**Last Class:**

$$Var(X+Y) = Var(X) + Var(Y)$$

· Linearity of variance.

· Markov's inequality: the most fundamental concentration bound. $\Pr(X \geq t \cdot \mathbb{E}[X]) \leq 1/t$.

· Algorithmic applications of Markov's inequality, linearity of expectation, and indicator random variables:

  · Counting collisions to estimate CAPTCHA database size.
  · Start on analyzing hash tables with random hash functions.
  · Collisions free hashing using a table with $O(m^2)$ slots to store $m$ items.

## Content Overview

**Today:**

- Finish up random hash functions and hash tables.
- 2-level hashing, 2-universal and pairwise independent hash functions.
- Application of random hashing to distributed load balancing.
- Through this application learn about Chebyshev's inequality, which strengthens Markov's inequality (maybe not until next class).

# Quiz Questions

The expected number of inches of rain on Saturday is 4 and the expected number of inches on Sunday is 2. There is a 50% chance of rain on Saturday. If it rains on Saturday, there is a 75% chance of rain on Sunday. If it does not rain on Saturday, there is only a 25% chance of rain on Sunday. What is the expected number of inches of rainfall total over the weekend?

Answer: [                ]

Check

$X =$ # inches on Sat

$Y =$ # inches on Sun

$$E[X+Y] = E[X] + E[Y] = 6$$
$$\qquad\quad\ \ 4 \quad + \ 2$$

# Hash Tables

We store *m* items from a large universe in a hash table with *n* positions.



- Want to show that when $h : U \rightarrow [n]$ is a fully random hash function, query time is $O(1)$ with good probability.
- Equivalently: want to show that there are few collisions between hashed items.

## Collision Free Hashing

Let $C = \sum_{i,j \in [m], i < j} C_{i,j}$ be the number of pairwise collisions between items.

$\mathbb{E}[C_{i,j}] = \frac{1}{n}$

$$\mathbb{E}[C] = \frac{m(m-1)}{2n} \quad \text{(via the Captcha analysis)}$$

- For $n = 4m^2$ we have: $\mathbb{E}[C] = \frac{m(m-1)}{8m^2} \leq \frac{1}{8}$.

$m$: total number of stored items, $n$: hash table size, $C$: total pairwise collisions in table.

## Collision Free Hashing

Let $C = \sum_{i,j \in [m], i < j} C_{i,j}$ be the number of pairwise collisions between items.

$$\mathbb{E}[C] = \frac{m(m-1)}{2n} \quad \text{(via the Captcha analysis)}$$

- For $n = 4m^2$ we have: $\mathbb{E}[C] = \frac{m(m-1)}{8m^2} \leq \frac{1}{8}$.

$x = 1$

Apply Markov's Inequality: $\Pr[C \geq 1] \leq \frac{\mathbb{E}[C]}{1} \leq \frac{1}{8}$.

> $m$: total number of stored items, $n$: hash table size, $C$: total pairwise collisions in table.

8

## Collision Free Hashing

Let $C = \sum_{i,j \in [m], i < j} C_{i,j}$ be the number of pairwise collisions between items.

$$\mathbb{E}[C] = \frac{m(m-1)}{2n} \quad \text{(via the Captcha analysis)}$$

- For $n = 4m^2$ we have: $\mathbb{E}[C] = \frac{m(m-1)}{8m^2} \leq \frac{1}{8}$.

**Apply Markov's Inequality:** $\Pr[C \geq 1] \leq \frac{\mathbb{E}[C]}{1} = \frac{1}{8}$.

$$\Pr[C = 0] = 1 - \Pr[C \geq 1] \geq 1 - \frac{1}{8} = \frac{7}{8}$$

$m$: total number of stored items, $n$: hash table size, $C$: total pairwise collisions in table.

## Collision Free Hashing

Let $C = \sum_{i,j \in [m], i < j} C_{i,j}$ be the number of pairwise collisions between items.

$$\mathbb{E}[C] = \frac{m(m-1)}{2n} \quad \text{(via the Captcha analysis)}$$

- For $n = 4m^2$ we have: $\mathbb{E}[C] = \frac{m(m-1)}{8m^2} \leq \frac{1}{8}$.

**Apply Markov's Inequality:** $\Pr[C \geq 1] \leq \frac{\mathbb{E}[C]}{1} = \frac{1}{8}$.

$$\Pr[C = 0] = 1 - \Pr[C \geq 1] \geq 1 - \frac{1}{8} = \frac{7}{8}$$

I.e., with probability at least 7/8 we have no collisions and thus $O(1)$ query time. But we are using $O(m^2)$ space to store $m$ items...

> $m$: total number of stored items, $n$: hash table size, $C$: total pairwise collisions in table.

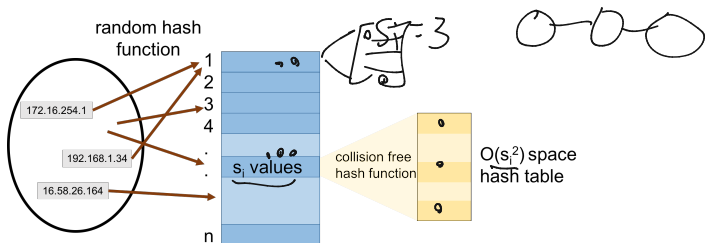## Two Level Hashing

Want to preserve $O(1)$ query time while using $O(m)$ space.

Want to preserve $O(1)$ query time while using $O(m)$ space.

## Two-Level Hashing:

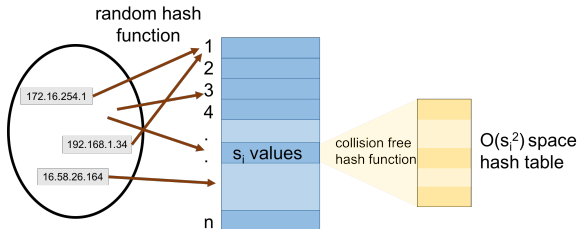Want to preserve $O(1)$ query time while using $O(m)$ space.

**Two-Level Hashing:**



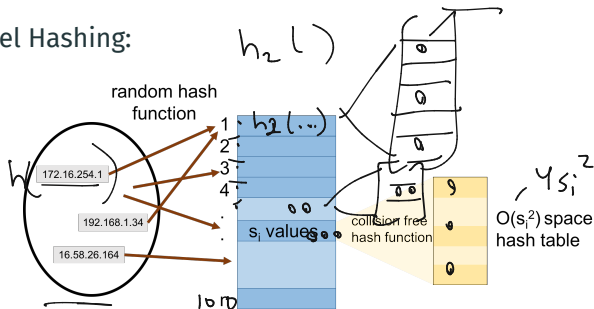- For each bucket with $s_i$ values, pick a collision free hash function mapping $[s_i] \to [s_i^2]$.

Want to preserve $O(1)$ query time while using $O(m)$ space.

Static
hash

**Two-Level Hashing:**

$h_2( )$



random hash
function

$h$

172.16.254.1

192.168.1.34

16.58.26.164

$h_2(...)$

$s_i$ values

collision free
hash function

$4s_i^2$

$O(s_i^2)$ space
hash table

$10m$

· For each bucket with $s_i$ values, pick a collision free hash
  function mapping $[s_i] \rightarrow [s_i^2]$.

· **Just Showed:** A random function is collision free with
  probability $\geq \frac{7}{8}$ so can just generate a <u>random hash function</u>
  and check if it is collision free.

Exercise: $O(\log n)$
hash functions

9

## Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions.

$x_j, x_k$: stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored in hash table at position $i$.
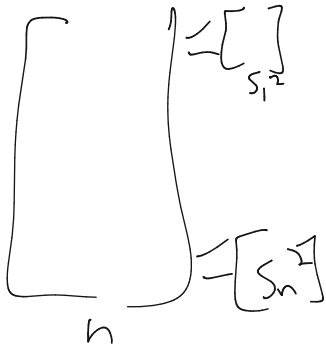
10

## Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

$x_j, x_k$: stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbf{S} = n + \sum_{i=1}^{n} \mathbf{s}_i^2$



$x_j, x_k$: stored items, $n$: hash table size, $\mathbf{h}$: random hash function, $\mathbf{S}$: space usage of two level hashing, $\mathbf{s}_i$: # items stored in hash table at position $i$.

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i] = \frac{m}{n}$$

$$\mathbb{E}[\mathsf{s}_i] = \sum_{j=1}^{m} \mathbb{1}_{i,j} = \sum_{j=1}^{m} \frac{1}{n} = \frac{m}{n}$$

$$\left( \frac{1}{n} \right)$$

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

## Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[S] = n + \sum_{i=1}^{n} \mathbb{E}[s_i^2]$

$x_j, x_k$: stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$\mathbb{I}_{\mathsf{h}(x_j)=i} = 1$ if item $j$ heshes into bucket $i$

$= 0$ o.w.

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[S] = n + \sum_{i=1}^{n} \mathbb{E}[s_i^2]$

$$\mathbb{E}[s_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{h(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k \in [m]} \mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right]$$

$\mathbb{I}_1 \cdot \mathbb{I}_1$

$(\mathbb{I}_1 + \mathbb{I}_2 + \mathbb{I}_3)(\mathbb{I}_1 + \mathbb{I}_2 + \mathbb{I}_3)$

$\mathbb{I}_1^2 + \mathbb{I}_1 \mathbb{I}_2 + \ldots \mathbb{I}_2 \cdot \mathbb{I}_1$

Collisions again!

---

$x_j, x_k$: stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\underline{\mathbb{E}[\mathsf{s}_i^2]} = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k\in[m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \underline{\mathbb{I}_{\mathsf{h}(x_k)=i}}\right].$$

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k \in [m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{\underbrace{j,k \in [m]}} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right].$$

· For $j = k$,

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. <span style="color:orange">What is the expected space usage?</span>

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k\in[m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{j,k\in[m]} \mathbb{E}\left[\underbrace{\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}}\right] \cdot$$

· For $j = k$,

$$\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \mathbb{E}\left[\left(\underbrace{\mathbb{I}_{\mathsf{h}(x_j)=i}}\right)^2\right]$$

$$= \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i}\right] = \frac{1}{n}$$

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k\in[m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right].$$

- For $j = k$,

$$\mathbb{E}\left[\underbrace{\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}}\right] = \mathbb{E}\left[\left(\mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right] = \underbrace{\Pr[\mathsf{h}(x_j) = i]}$$

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k\in[m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right].$$

- For $j = k$,

$$\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \mathbb{E}\left[\left(\mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right] = \Pr[\mathsf{h}(x_j) = i] = \frac{1}{n}.$$

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k\in[m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right].$$

- For $j = k$,

  $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \mathbb{E}\left[\left(\mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right] = \Pr[\mathsf{h}(x_j) = i] = \frac{1}{n}.$

- For $j \neq k$,

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[ \left( \sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i} \right)^2 \right]$$

$$= \mathbb{E}\left[ \sum_{j,k \in [m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i} \right] = \sum_{j,k \in [m]} \mathbb{E}\left[ \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i} \right].$$

- For $j = k$,
  $\mathbb{E}\left[ \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i} \right] = \mathbb{E}\left[ \left( \mathbb{I}_{\mathsf{h}(x_j)=i} \right)^2 \right] = \Pr[\mathsf{h}(x_j) = i] = \frac{1}{n}$.
- For $j \neq k$, $\mathbb{E}\left[ \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i} \right]$

  $0 \cdot 1 \quad 0 \cdot 1 \quad = 1$ if $x_j$ & $x_k$ hash to $i$

  $0 \quad$ o.w.

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

# Space Usage

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$$\mathbb{E}[\mathsf{s}_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right]$$

$$= \mathbb{E}\left[\sum_{j,k \in [m]} \mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right].$$

- For $j = k$,
  $$\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \mathbb{E}\left[\left(\mathbb{I}_{\mathsf{h}(x_j)=i}\right)^2\right] = \Pr[\mathsf{h}(x_j) = i] = \frac{1}{n}.$$

- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \underbrace{\Pr[\mathsf{h}(x_j) = i} \cap \underbrace{\mathsf{h}(x_k) = i]}$

---

$x_j, x_k$: stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored in hash table at position $i$.

Query time for two level hashing is $O(1)$: requires evaluating two hash functions. What is the expected space usage?

Up to constants, space used is: $\mathbb{E}[S] = n + \sum_{i=1}^{n} \mathbb{E}[s_i^2]$

$$\mathbb{E}[s_i^2] = \mathbb{E}\left[\left(\sum_{j=1}^{m} \mathbb{I}_{h(x_j)=i}\right)^2\right]$$

*(handwritten annotations:)*
$= \mathbb{E}\left[\left(\mathbb{I}_1 + \mathbb{I}_2 + \mathbb{I}_3\right)^2\right]$
$\mathbb{E}\left[\mathbb{I}_1^2 + \mathbb{I}_2^2 + \mathbb{I}_3^2 + 2\mathbb{I}_1\mathbb{I}_2 + 2\mathbb{I}_2\mathbb{I}_3 + 2\mathbb{I}_1\mathbb{I}_3\right]$

$$= \mathbb{E}\left[\sum_{j,k\in[m]} \mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right].$$

- For $j = k$,
  $$\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \mathbb{E}\left[\left(\mathbb{I}_{h(x_j)=i}\right)^2\right] = \Pr[h(x_j) = i] = \tfrac{1}{n}.$$

- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \Pr[h(x_j) = i \cap h(x_k) = i] = \tfrac{1}{n^2}.$

$x_j, x_k$: stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored in hash table at position $i$.

# Space Usage

$$\mathbb{E}[\mathsf{s}_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right]$$

$$\sum_{j=1}^{m} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i}^{2}\right] + \sum_{\substack{j,k \in [m] \\ j \neq k}} \mathbb{E}\left[\mathbb{I}_j \cdot \mathbb{I}_k\right]$$

$$\frac{1}{n} \qquad \frac{1}{n^2}$$

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n}$.
- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n^2}$.

---

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored at pos $i$.

## Space Usage

$$\mathbb{E}[\mathsf{s}_i^2] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right]$$

$$= \underbrace{m \cdot \frac{1}{n}} + \underbrace{2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}}$$

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n}$.
- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n^2}$.

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored at pos $i$.

# Space Usage

$$\mathbb{E}[s_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right]$$

$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n}$.

- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n^2}$.

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored at pos $i$.

# Space Usage

$$\mathbb{E}[s_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right]$$

$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n}$.

- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n^2}$.

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored at pos $i$.

# Space Usage

$$\mathbb{E}[\mathsf{s}_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right]$$

$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

$$= \frac{m}{n} + \frac{m(m-1)}{n^2} \quad +n$$

*expected* *space'*

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n}$.
- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n^2}$.

---

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored at pos $i$.

## Space Usage

$$\mathbb{E}[\mathsf{s}_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right]$$

$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

$$= \underbrace{\frac{m}{n}}_{} + \underbrace{\frac{m(m-1)}{n^2}}_{} \leq 2 \text{ (If we set } n = m.)$$

· For $j = k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n}$.

· For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right] = \frac{1}{n^2}$.

---

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $\mathsf{h}$: random hash function, $\mathsf{S}$: space usage of two level hashing, $\mathsf{s}_i$: # items stored at pos $i$.

# Space Usage

$$\mathbb{E}[s_i^2] = \sum_{j,k\in[m]} \mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right]$$

$$/) \qquad = m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

$$2 \qquad = \frac{m}{n} + \frac{m(m-1)}{n^2} \leq 2 \text{ (If we set } n = m.)$$

· For $j = k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n}$.

· For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n^2}$.

Total Expected Space Usage: (if we set $n = m$)

$$\underbrace{\mathbb{E}[S]} = \underbrace{n} + \sum_{i=1}^{n} \overset{2}{\mathbb{E}[s_i^2]}$$

---

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, $h$: random hash function, $S$: space usage of two level hashing, $s_i$: # items stored at pos $i$.

# Space Usage

$$\mathbb{E}[s_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right]$$

$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

$$\frac{m}{n}^1 + \frac{m(m-1)}{m^2}^{\leq 1}$$

$$\mathbb{E}[s_i^2] = \frac{m}{n} + \frac{m(m-1)}{n^2} \leq 2 \text{ (If we set } n = m.\text{)}$$

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n}$.
- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n^2}$.

Total Expected Space Usage: (if we set $n = m$)

$$\mathbb{E}[S] = n + \sum_{i=1}^{n} \mathbb{E}[s_i^2] \leq n + n \cdot 2 = 3n = 3m.$$

$x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, h: random hash function, S: space usage of two level hashing, $s_i$: # items stored at pos $i$.

# Space Usage

$$\mathbb{E}[s_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right]$$

$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

$$= \frac{m}{n} + \frac{m(m-1)}{n^2} \leq 2 \text{ (If we set } n = m.)$$

- For $j = k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n}$.
- For $j \neq k$, $\mathbb{E}\left[\mathbb{I}_{h(x_j)=i} \cdot \mathbb{I}_{h(x_k)=i}\right] = \frac{1}{n^2}$.

**Total Expected Space Usage:** (if we set $n = m$)

$$\mathbb{E}[S] = n + \sum_{i=1}^{n} \mathbb{E}[s_i^2] \leq n + n \cdot 2 = 3n = 3m.$$

Near optimal space with $O(1)$ query time!

> $x_j, x_k$: stored items, $m$: # stored items, $n$: hash table size, h: random hash function, S: space usage of two level hashing, $s_i$: # items stored at pos $i$.

## Efficiently Computable Hash Function

So Far: we have assumed a fully random hash function $h(x)$ with $\Pr[h(x) = i] = \frac{1}{n}$ for $i \in 1, \ldots, n$ and $h(x), h(y)$ independent for $x \neq y$.

## Efficiently Computable Hash Function

So Far: we have assumed a **fully random hash function** $h(x)$ with $\Pr[h(x) = i] = \frac{1}{n}$ for $i \in 1, \ldots, n$ and $h(x), h(y)$ independent for $x \neq y$.

- To compute a random hash function we have to store a table of $x$ values and their hash values. Would take at least $O(m)$ space and $O(m)$ query time to look up $h(x)$ if we hash $m$ values. Making our whole quest for $O(1)$ query time pointless!

| x | h(x) |
|---|---|
| $x_1$ | 45 |
| $x_2$ | 1004 |
| $x_3$ | 10 |
| ⋮ | ⋮ |
| $x_m$ | 12 |

## Efficiently Computable Hash Functions

What properties did we use of the randomly chosen hash function?

## Efficiently Computable Hash Functions

What properties did we use of the randomly chosen hash function?

> **Pairwise Independent Hash Function.** A random hash function
> from $h : U \to [n]$ is pairwise independent if for all $i, j \in [n]$:
>
> $$\Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}.$$

## Efficiently Computable Hash Functions

What properties did we use of the randomly chosen hash function?

**Pairwise Independent Hash Function.** A random hash function from $h : U \to [n]$ is pairwise independent if for all $i, j \in [n]$:

$$\Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}.$$

Exercise 1: Check the two-level hashing proof to confirm that this property is all that was needed.

What properties did we use of the randomly chosen hash function?

> **Pairwise Independent Hash Function.** A random hash function
> from $h : U \to [n]$ is pairwise independent if for all $i, j \in [n]$:
>
> $$x \neq y \qquad \Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}.$$

**Exercise 1:** Check the two-level hashing proof to confirm that this
property is all that was needed.

When $h(x)$ and $h(y)$ are chosen independently at random from $[n]$,
$\Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}$ (so a fully random hash function is
pairwise independent).

# Efficiently Computable Hash Functions

What properties did we use of the randomly chosen hash function?

> **Pairwise Independent Hash Function.** A random hash function
> from $h : U \to [n]$ is pairwise independent if for all $i, j \in [n]$:
>
> $$\Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}.$$

**Exercise 1:** Check the two-level hashing proof to confirm that this property is all that was needed.

When $h(x)$ and $h(y)$ are chosen independently at random from $[n]$, $\Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}$ (so a fully random hash function is pairwise independent).

**Efficient Implementation:** Let $p$ be a prime with $p \geq |U|$. Choose random $a, b \in [p]$ with $a \neq 0$. Represent $x$ as an integer and let

$$h(x) = (ax + b \mod p) \mod n.$$

Another common requirement for a hash function:

> **2-Universal Hash Function** (low collision probability). A random hash function from $h : U \rightarrow [n]$ is two universal if:
>
> $x \neq y$
>
> $$\Pr[h(x) = h(y)] \leq \frac{1}{n}.$$

## Universal Hashing

Another common requirement for a hash function:

> **2-Universal Hash Function** (low collision probability). A random hash function from $h : U \to [n]$ is two universal if:
>
> $$\Pr[h(x) = h(y)] \le \frac{1}{n}.$$

Think-Pair-Shair: Which is a more stringent requirement?
2-universal or pairwise independent?

## Universal Hashing

Another common requirement for a hash function:

> **2-Universal Hash Function** (low collision probability). A random hash function from $h : U \to [n]$ is two universal if:
>
> $$\Pr[h(x) = h(y)] \leq \frac{1}{n}.$$

**Think-Pair-Shair:** Which is a more stringent requirement? 2-universal or pairwise independent?

> **Pairwise Independent Hash Function.** A random hash function from $h : U \to [n]$ is pairwise independent if for all $i, j \in [n]$:
>
> $$\Pr[h(x) = i \cap h(y) = j] = \frac{1}{n^2}.$$

## Universal Hashing

Another common requirement for a hash function:

> **2-Universal Hash Function** (low collision probability). A random hash function from $h : U \to [n]$ is two universal if:
>
> $$Pr[h(x) = h(y)] \leq \frac{1}{n}.$$

**Think-Pair-Shair:** Which is a more stringent requirement?
2-universal or pairwise independent?

$$Pr[h(x) = h(y)] = \sum_{i=1}^{n} Pr[h(x) = i \cap h(y) = i] = n \cdot \frac{1}{n^2} = \frac{1}{n}.$$

## Universal Hashing

Another common requirement for a hash function:

> **2-Universal Hash Function** (low collision probability). A random hash function from $h : U \to [n]$ is two universal if:
> $$\Pr[h(x) = h(y)] \leq \frac{1}{n}.$$

**Think-Pair-Shair:** Which is a more stringent requirement?
2-universal or pairwise independent?

$$\Pr[h(x) = h(y)] = \sum_{i=1}^{n} \Pr[h(x) = i \cap h(y) = i] = n \cdot \frac{1}{n^2} = \frac{1}{n}.$$

**Remember:** A fully random hash function is both 2-universal and pairwise independent. But it is not efficiently implementable.

Another common requirement for a hash function:

> **2-Universal Hash Function** (low collision probability). A random hash function from $h : U \to [n]$ is two universal if:
>
> $$\Pr[h(x) = h(y)] \leq \frac{1}{n}.$$

**Think-Pair-Shair:** Which is a more stringent requirement?
2-universal or pairwise independent?

$$\Pr[h(x) = h(y)] = \sum_{i=1}^{n} \Pr[h(x) = i \cap h(y) = i] = n \cdot \frac{1}{n^2} = \frac{1}{n}.$$

**Remember:** A fully random hash function is both 2-universal and pairwise independent. But it is not efficiently implementable.

**Exercise 2:** Rework the two-level hashing proof to show that 2-universality is in fact all that is needed.