## COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2020.

Lecture 3

By Thursday:

- Sign up for Piazza.
- Sign up for Gradescope (code on class website) and fill out the Gradescope consent poll on Piazza. Contact me via email if you don't consent to use Gradescope.

### By Thursday:

- Sign up for Piazza.
- Sign up for Gradescope (code on class website) and fill out the Gradescope consent poll on Piazza. Contact me via email if you don't consent to use Gradescope.

First Problem Set: released Saturday, due 9/11 at 8pm in Gradescope.

- Remember you can complete in a group of up to 3 students, who all turn in one submission with three names on it.

91 students completed the quizzes – make sure that if you are enrolled you are doing the quiz each week.

91 students completed the quizzes – make sure that if you are enrolled you are doing the quiz each week.

**Question 1:** The expected number of inches of rain on Saturday is 2 and the expected number of inches on Sunday is 6. 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?

91 students completed the quizzes – make sure that if you are enrolled you are doing the quiz each week.

**Question 1:** The expected number of inches of rain on Saturday is 2 and the expected number of inches on Sunday is 6. 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?

**Concerns:** Probability/linear algebra background, proofs/derivations.

Last Class We Covered:

- Markov's inequality: the most fundamental concentration bound.
- Algorithmic applications of Markov's inequality, linearity of expectation, and indicator random variables:
  - Counting collisions to estimate CAPTCHA database size.
  - Counting collisions to understand the runtime of hash tables with random hash functions.
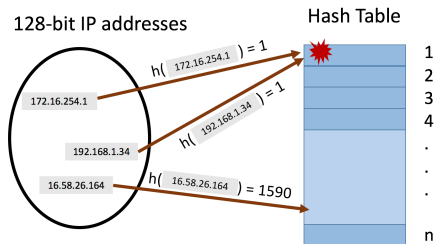
Last Class We Covered:

- Markov's inequality: the most fundamental concentration bound.
- Algorithmic applications of Markov's inequality, linearity of expectation, and indicator random variables:
  - Counting collisions to estimate CAPTCHA database size.
  - Counting collisions to understand the runtime of hash tables with random hash functions.
- Collision counting is closely related to the birthday paradox.

Today:

- Finish up random hash functions and hash tables.
- See an applications of random hashing to load balancing in distributed systems.
- Through these applications learn about:
  - Chebyshev's inequality, which strengthens Markov's inequality.
  - The union bound, for understanding the probabilities of correlated random events.

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



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

When storing *m* items in a table of size *n*, the expected number of pairwise collisions (two items stored in the same slots) is:

$$\mathbb{E}[\mathsf{C}] = \frac{m(m-1)}{2n}.$$

*m*: total number of stored items, *n*: hash table size, $\mathsf{C}$: total pairwise collisions in table.

When storing $m$ items in a table of size $n$, the expected number of pairwise collisions (two items stored in the same slots) is:

$$\mathbb{E}[C] = \frac{m(m-1)}{2n}.$$

- For $n = 4m^2$ we have: $\mathbb{E}[C] = \frac{m(m-1)}{8m^2} \leq \frac{1}{8}$.
- By Markov's inequality there no collisions with probability at least $\frac{7}{8}$.

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

When storing *m* items in a table of size *n*, the expected number of pairwise collisions (two items stored in the same slots) is:

$$\mathbb{E}[\mathsf{C}] = \frac{m(m-1)}{2n}.$$

- For $n = 4m^2$ we have: $\mathbb{E}[\mathsf{C}] = \frac{m(m-1)}{8m^2} \leq \frac{1}{8}$.
- By Markov's inequality there no collisions with probability at least $\frac{7}{8}$.

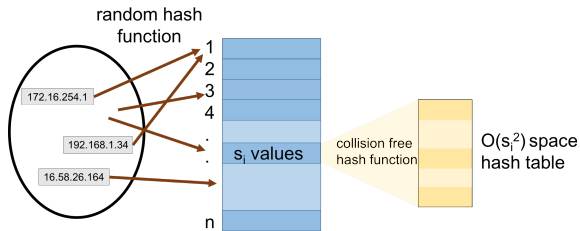$O(1)$ query time, but we are using $O(m^2)$ space to store *m* items...

> *m*: total number of stored items, *n*: hash table size, $\mathsf{C}$: total pairwise collisions in table.

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:
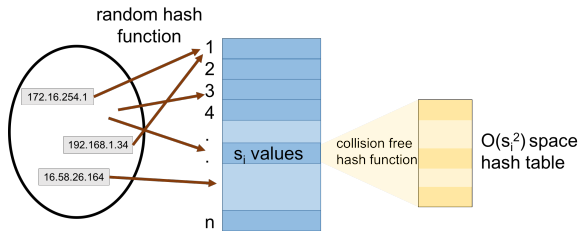


random hash function

1
2
3
4
.
.

172.16.254.1

192.168.1.34

16.58.26.164

n

$s_i$ values

collision free hash function

$O(s_i^2)$ space hash table

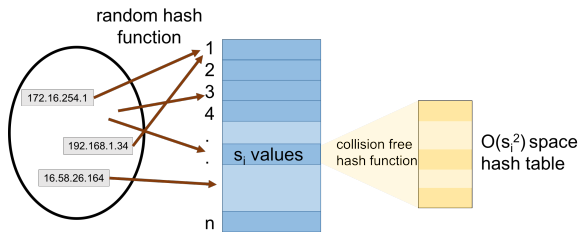Want to preserve $O(1)$ query time while using $O(m)$ space.

**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.

## Two-Level Hashing:



- For each bucket with $s_i$ values, pick a collision free hash function mapping $[s_i] \to [s_i^2]$.
- **Just Showed:** A random function is collision free with probability $\geq \frac{7}{8}$ so can just generate a random hash function and check if it is collision free.

7

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 $i$.

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 $i$.

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: $\mathsf{S} = n + \sum_{i=1}^{n} \mathsf{s}_i^2$

$x_j, x_k$: 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: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$x_j, x_k$: 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: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

$x_j, x_k$: 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: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

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

$x_j, x_k$: 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: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

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

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

Collisions again!

$x_j, x_k$: 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: $\mathbb{E}[\mathsf{S}] = n + \sum_{i=1}^{n} \mathbb{E}[\mathsf{s}_i^2]$

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

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

$x_j, x_k$: 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 $j = k$,

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

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

$x_j, x_k$: 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 $i$.

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

$x_j, x_k$: 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: $\mathbb{E}[S] = n + \sum_{i=1}^{n} \mathbb{E}[s_i^2]$

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

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

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

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

---

$x_j, x_k$: 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 $i$.

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

$x_j, x_k$: 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 $i$.

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

- 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 $i$.

9

$$\mathbb{E}[\mathsf{s}_i^2] = \sum_{j,k \in [m]} \mathbb{E}\left[\mathbb{I}_{\mathsf{h}(x_j)=i} \cdot \mathbb{I}_{\mathsf{h}(x_k)=i}\right]$$
$$= m \cdot \frac{1}{n} + 2 \cdot \binom{m}{2} \cdot \frac{1}{n^2}$$

· 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 $i$.

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

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

- 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 $i$.

9

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

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

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

---

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

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

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

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

· 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 $i$.

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

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

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

- 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 $i$.

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

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

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

- 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 $n = m$)

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

$x_j, x_k$: 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 $i$.

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

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

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

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

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

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

---

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

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

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

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

- 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 $i$.

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 $x \neq y$.

So Far: we have assumed a **fully random hash function** $h(x)$ with $\Pr[h(x) = i] = \frac{1}{n}$ for $i \in 1, \dots, n$ 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 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 |
| $\vdots$ | $\vdots$ |
| $x_m$ | 12 |

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

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

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

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

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

**Exercise:** Rework the two level hashing proof to show that this property is really all that is needed.

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

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

**Exercise:** Rework the two level hashing proof to show that this property is really all that is needed.

When $h(x)$ and $h(y)$ are chosen independently at random from $[n]$, $\Pr[h(x) = h(y)] = \frac{1}{n}$ (so a fully random hash function is 2-universal)

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

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

**Exercise:** Rework the two level hashing proof to show that this property is really all that is needed.

When $h(x)$ and $h(y)$ are chosen independently at random from $[n]$, $\Pr[h(x) = h(y)] = \frac{1}{n}$ (so a fully random hash function is 2-universal)

**Efficient Alternative:** Let $p$ be a prime with $p \geq |U|$. Choose random $a, b \in [p]$ with $a \neq 0$. Let:

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

11

Another common requirement for a hash function:

Another common requirement for a hash function:

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

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

Another common requirement for a hash function:

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

**Breakout:** Which is a more stringent requirement? 2-universal or pairwise independent?

Another common requirement for a hash function:

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

**Breakout:** Which is a more stringent requirement? 2-universal or
pairwise independent?

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

Another common requirement for a hash function:

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

**Breakout:** Which is a more stringent requirement? 2-universal or pairwise independent?

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

A closely related $(\mathsf{a}x + \mathsf{b}) \mod p$ construction gives pairwise independence on top of 2-universality.

Another common requirement for a hash function:

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

**Breakout:** Which is a more stringent requirement? 2-universal or
pairwise independent?

$$\Pr[h(x) = h(y)] = \sum_{i=1}^{n} \Pr[h(x) = i \cap h(y) = i] = n \cdot \frac{1}{n^2} = \frac{1}{n}.$$
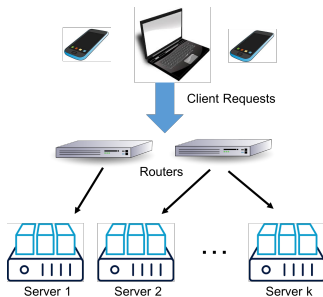
A closely related $(ax + b) \mod p$ construction gives pairwise
independence on top of 2-universality.

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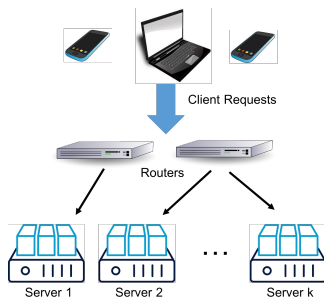pairwise independent. But it is not efficiently implementable.

1. We'll consider an application where our toolkit of linearity of expectation + Markov's inequality doesn't give much.

1. We'll consider an application where our toolkit of linearity of expectation + Markov's inequality doesn't give much.

2. Then we'll show how a simple twist on Markov's can give a much stronger result.
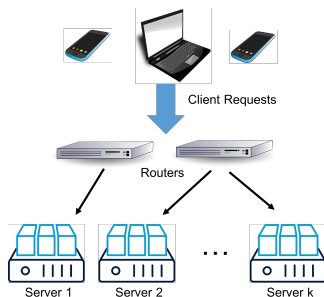
Randomized Load Balancing:



Client Requests

Routers

Server 1    Server 2    . . .    Server k

Randomized Load Balancing:



Client Requests

Routers

Server 1    Server 2    ...    Server k

**Simple Model:** $n$ requests randomly assigned to $k$ servers. How many requests must each server handle?

Randomized Load Balancing:



Client Requests

Routers

Server 1   Server 2   ...   Server k

**Simple Model:** $n$ requests randomly assigned to $k$ servers. How many requests must each server handle?

· Often assignment is done via a random hash function. Why?

Expected Number of requests assigned to server $i$:

$$\mathbb{E}[\mathsf{R}_i] = \sum_{j=1}^{n} \mathbb{E}[\mathbb{I}_{\text{request } j \text{ assigned to } i}] = \sum_{j=1}^{n} \Pr\left[j \text{ assigned to } i\right] = \frac{n}{k}.$$

$n$: total number of requests, $k$: number of servers randomly assigned requests, $\mathsf{R}_i$: number of requests assigned to server $i$.

Expected Number of requests assigned to server $i$:

$$\mathbb{E}[\mathsf{R}_i] = \sum_{j=1}^{n} \mathbb{E}[\mathbb{I}_{\text{request } j \text{ assigned to } i}] = \sum_{j=1}^{n} \Pr[j \text{ assigned to } i] = \frac{n}{k}.$$

If we provision each server be able to handle twice the expected load, what is the probability that a server is overloaded?

$n$: total number of requests, $k$: number of servers randomly assigned requests, $\mathsf{R}_i$: number of requests assigned to server $i$.

Expected Number of requests assigned to server $i$:

$$\mathbb{E}[\mathsf{R}_i] = \sum_{j=1}^{n} \mathbb{E}[\mathbb{I}_{\text{request } j \text{ assigned to } i}] = \sum_{j=1}^{n} \Pr[j \text{ assigned to } i] = \frac{n}{k}.$$

If we provision each server be able to handle twice the expected load, what is the probability that a server is overloaded?

Applying Markov's Inequality

$$\Pr[\mathsf{R}_i \geq 2\mathbb{E}[\mathsf{R}_i]] \leq \frac{\mathbb{E}[\mathsf{R}_i]}{2\mathbb{E}[\mathsf{R}_i]} = \frac{1}{2}.$$

> $n$: total number of requests, $k$: number of servers randomly assigned requests,
> $\mathsf{R}_i$: number of requests assigned to server $i$.

Expected Number of requests assigned to server $i$:

$$\mathbb{E}[\mathsf{R}_i] = \sum_{j=1}^{n} \mathbb{E}[\mathbb{I}_{\text{request } j \text{ assigned to } i}] = \sum_{j=1}^{n} \Pr[j \text{ assigned to } i] = \frac{n}{k}.$$

If we provision each server be able to handle twice the expected load, what is the probability that a server is overloaded?

Applying Markov's Inequality

$$\Pr[\mathsf{R}_i \geq 2\mathbb{E}[\mathsf{R}_i]] \leq \frac{\mathbb{E}[\mathsf{R}_i]}{2\mathbb{E}[\mathsf{R}_i]} = \frac{1}{2}.$$

Not great...half the servers may be overloaded.

> $n$: total number of requests, $k$: number of servers randomly assigned requests, $\mathsf{R}_i$: number of requests assigned to server $i$.