

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Spring 2026.

Lecture 4

- Problem Set 1 due next Friday 2/20, at 11:59pm.
- For the Wikipedia Mark-and-Recapture question I posted a file with random URLs in case you are having trouble with requests to Wikipedia:Random.
- I am out of town next week. So Tuesday's lecture will be over Zoom. Link will be posted in Piazza. No lecture Thursday (Monday schedule for UMass).
- I will hold my usual Tuesday 2:30pm office hours, over the same Zoom link.

Last Time

Last Class:

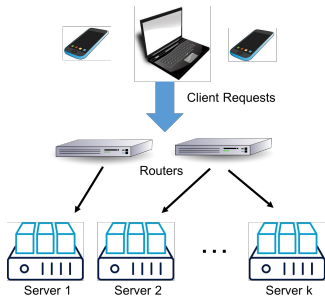
- 2-level hashing and its analysis via linearity of expectation.
Gives optimal $O(1)$ query time and $O(m)$ expected space usage.
- Practical random hash functions: 2-universal and pairwise independent hashing.

This Time:

- Hashing for load balancing in distributed systems. Motivating:
 - Stronger concentration inequalities: Chebyshev's inequality, exponential tail bounds, and their connections to the law of **large numbers and central limit theorem**.
 - The union bound to bound the probability that one of multiple possible correlated events happens.
- Some of the problem set questions use Chebyshev's inequality. After today you will be able to solve them.

Another Application

Randomized Load Balancing:



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? Why not just a random number generator?

Weakness of Markov's

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

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

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

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

n : total number of requests, k : number of servers randomly assigned requests,
 R_i : number of requests assigned to server i .

Chebyshev's inequality

With a very simple twist, Markov's inequality can be made much more powerful.

For any random variable X and any value $t > 0$:

$$\Pr(|X| \geq t) = \Pr(X^2 \geq t^2).$$

X^2 is a nonnegative random variable. So can apply Markov's inequality:

Chebyshev's inequality:

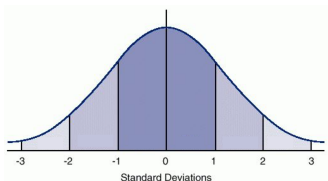
$$\Pr(|X - \mathbb{E}[X]| \geq t) = \Pr(X^2 \geq t^2) \leq \frac{\mathbb{E}[X^2]}{t^2} = \frac{\text{Var}[X]}{t^2}.$$

(by plugging in the random variable $X - \mathbb{E}[X]$)

Chebyshev's inequality

$$\Pr(|X - \mathbb{E}[X]| \geq t) \leq \frac{\text{Var}[X]}{t^2}$$

What is the probability that X falls s standard deviations from its mean?



$$\Pr(|X - \mathbb{E}[X]| \geq s \cdot \sqrt{\text{Var}[X]}) \leq \frac{\text{Var}[X]}{s^2 \cdot \text{Var}[X]} = \frac{1}{s^2}.$$

X : any random variable, t, s : any fixed numbers.

Law of Large Numbers

Consider drawing independent identically distributed (i.i.d.) random variables X_1, \dots, X_n with mean μ and variance σ^2 .

How well does the sample average $S = \frac{1}{n} \sum_{i=1}^n X_i$ approximate the true mean μ ?

$$\text{Var}[S] = \text{Var} \left[\frac{1}{n} \sum_{i=1}^n X_i \right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{1}{n^2} \cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}.$$

By Chebyshev's Inequality: for any fixed value $\epsilon > 0$,

$$\Pr(|S - \mathbb{E}[S]| \geq \epsilon) \leq \frac{\text{Var}[S]}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}.$$

Law of Large Numbers: with enough samples n , the sample average will always concentrate to the mean.

- Cannot show from vanilla Markov's inequality.

Load Balancing Variance

We can write the number of requests assigned to server i , R_i as:

$$R_i = \sum_{j=1}^n R_{i,j} \quad \text{Var}[R_i] = \sum_{j=1}^n \text{Var}[R_{i,j}] \quad (\text{linearity of variance})$$

where $R_{i,j}$ is 1 if request j is assigned to server i and 0 otherwise.

$$\begin{aligned} \text{Var}[R_{i,j}] &= \mathbb{E}[R_{i,j}^2] - \mathbb{E}[R_{i,j}]^2 \\ &= \mathbb{E}[R_{i,j}] - \mathbb{E}[R_{i,j}]^2 \\ &= \frac{1}{k} - \frac{1}{k^2} \leq \frac{1}{k} \implies \text{Var}[R_i] \leq \frac{n}{k}. \end{aligned}$$

n : total number of requests, k : number of servers randomly assigned requests,
 R_i : number of requests assigned to server i .

Bounding the Load via Chebyshev's

Letting R_i be the number of requests sent to server i , $\mathbb{E}[R_i] = \frac{n}{k}$ and $\text{Var}[R_i] \leq \frac{n}{k}$.

Applying Chebyshev's:

$$\Pr\left(R_i \geq \frac{2n}{k}\right) \leq \Pr\left(|R_i - \mathbb{E}[R_i]| \geq \frac{n}{k}\right) \leq \frac{n/k}{n^2/k^2} = \frac{k}{n}.$$

- Overload probability is small when $k \ll n$!
- Might seem counterintuitive – bound gets worse as k grows.
- When k is large, the number of requests each server sees in expectation is very small so the law of large numbers doesn't 'kick in'.

n : total number of requests, k : number of servers randomly assigned requests,
 R_i : number of requests assigned to server i .

Maximum Server Load

What is the probability that the **maximum server load** exceeds $2 \cdot \mathbb{E}[\mathbf{R}_i] = \frac{2n}{k}$. I.e., that some server is overloaded if we give each $\frac{2n}{k}$ capacity?

$$\Pr \left(\max_i (\mathbf{R}_i) \geq \frac{2n}{k} \right) = \Pr \left(\left[\mathbf{R}_1 \geq \frac{2n}{k} \right] \cup \left[\mathbf{R}_2 \geq \frac{2n}{k} \right] \cup \dots \cup \left[\mathbf{R}_k \geq \frac{2n}{k} \right] \right) = \Pr$$

We want to show that $\Pr \left(\bigcup_{i=1}^k \left[\mathbf{R}_i \geq \frac{2n}{k} \right] \right)$ is small.

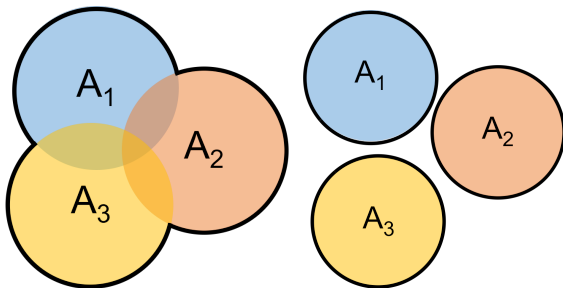
How do we do this? Note that $\mathbf{R}_1, \dots, \mathbf{R}_k$ are correlated in a somewhat complex way.

n : total number of requests, k : number of servers randomly assigned requests,
 \mathbf{R}_i : number of requests assigned to server i . $\mathbb{E}[\mathbf{R}_i] = \frac{n}{k}$. $\text{Var}[\mathbf{R}_i] = \frac{n}{k}$.

The Union Bound

Union Bound: For any random events A_1, A_2, \dots, A_k ,

$$\Pr(A_1 \cup A_2 \cup \dots \cup A_k) \leq \Pr(A_1) + \Pr(A_2) + \dots + \Pr(A_k).$$



When is the union bound tight? When A_1, \dots, A_k are all disjoint.

Applying the Union Bound

What is the probability that the **maximum server load** exceeds $2 \cdot \mathbb{E}[R_i] = \frac{2n}{k}$. I.e., that some server is overloaded if we give each $\frac{2n}{k}$ capacity?

$$\begin{aligned}\Pr\left(\max_i(R_i) \geq \frac{2n}{k}\right) &= \Pr\left(\bigcup_{i=1}^k \left[R_i \geq \frac{2n}{k}\right]\right) \\ &\leq \sum_{i=1}^k \Pr\left(\left[R_i \geq \frac{2n}{k}\right]\right) \quad (\text{Union Bound}) \\ &\leq \sum_{i=1}^k \frac{k}{n} = \frac{k^2}{n} \quad (\text{Bound from Chebyshev's})\end{aligned}$$

As long as $k \leq O(\sqrt{n})$, with good probability, the maximum server load will be small (compared to the expected load).

n : total number of requests, k : number of servers randomly assigned requests,
 R_i : number of requests assigned to server i . $\mathbb{E}[R_i] = \frac{n}{k}$. $\text{Var}[R_i] = \frac{n}{k}$.

Questions on union bound, Chebyshev's inequality,
random hashing?

Flipping Coins

We flip $n = 100$ independent coins, each are heads with probability $1/2$ and tails with probability $1/2$. Let H be the number of heads.

$$\mathbb{E}[H] = \frac{n}{2} = 50 \text{ and } \text{Var}[H] = \frac{n}{4} = 25 \rightarrow s.d. = 5$$

Markov's:	Chebyshev's:	In Reality:
$\Pr(H \geq 60) \leq .833$	$\Pr(H \geq 60) \leq .25$	$\Pr(H \geq 60) = 0.0284$
$\Pr(H \geq 70) \leq .714$	$\Pr(H \geq 70) \leq .0625$	$\Pr(H \geq 70) = .000039$
$\Pr(H \geq 80) \leq .625$	$\Pr(H \geq 80) \leq .0278$	$\Pr(H \geq 80) < 10^{-9}$

H has a simple Binomial distribution, so can compute these probabilities exactly.

Tighter Concentration Bounds

To be fair... Markov and Chebyshev's inequalities apply much more generally than to Binomial random variables like coin flips.

Can we obtain tighter concentration bounds that still apply to very general distributions?

- Markov's: $\Pr(X \geq t) \leq \frac{\mathbb{E}[X]}{t}$. **First Moment.**
- Chebyshev's: $\Pr(|X - \mathbb{E}[X]| \geq t) = \Pr(|X - \mathbb{E}[X]|^2 \geq t^2) \leq \frac{\text{Var}[X]}{t^2}$. **Second Moment.**
- What if we just apply Markov's inequality to even higher moments?

A Fourth Moment Bound

$$\Pr(|X - \mathbb{E}[X]| \geq t) = \Pr\left((X - \mathbb{E}[X])^4 \geq t^4\right) \leq \frac{\mathbb{E}\left[(X - \mathbb{E}[X])^4\right]}{t^4}.$$

Application to Coin Flips: Recall: $n = 100$ independent fair coins, \mathbf{H} is the number of heads.

- Bound the fourth moment:

$$\mathbb{E}\left[(\mathbf{H} - \mathbb{E}[\mathbf{H}])^4\right] = \mathbb{E}\left[\left(\sum_{i=1}^{100} \mathbf{H}_i - 50\right)^4\right] = \sum_{i,j,k,\ell} c_{ijkl} \mathbb{E}[\mathbf{H}_i \mathbf{H}_j \mathbf{H}_k \mathbf{H}_\ell] = 1862.5$$

where $\mathbf{H}_i = 1$ if coin flip i is heads and 0 otherwise. Then apply some messy calculations...

- Apply Fourth Moment Bound: $\Pr(|\mathbf{H} - \mathbb{E}[\mathbf{H}]| \geq t) \leq \frac{1862.5}{t^4}$.

Tighter Bounds

Chebyshev's:	4 th Moment:	In Reality:
$\Pr(H \geq 60) \leq .25$	$\Pr(H \geq 60) \leq .186$	$\Pr(H \geq 60) = 0.0284$
$\Pr(H \geq 70) \leq .0625$	$\Pr(H \geq 70) \leq .0116$	$\Pr(H \geq 70) = .000039$
$\Pr(H \geq 80) \leq .04$	$\Pr(H \geq 80) \leq .0023$	$\Pr(H \geq 80) < 10^{-9}$

Can we just keep applying Markov's inequality to higher and higher moments and getting tighter bounds?

- Yes! To a point.
- In fact – don't need to just apply Markov's to $|X - \mathbb{E}[X]|^k$ for some k . Can apply to any monotonic function $f(|X - \mathbb{E}[X]|)$.
- **Why monotonic?**

$$\Pr(|X - \mathbb{E}[X]| > t) = \Pr(f(|X - \mathbb{E}[X]|) > f(t)).$$

H: total number heads in 100 random coin flips. $\mathbb{E}[H] = 50$.

Next time: Use this approach to give exponential tail bounds, and a quantitative understanding of the central limit theorem.

- Leads to much better bounds for random hash tables, randomized load balancing, and many other randomized algorithms.