

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2021.

Lecture 4

- Problem set one posted this morning. Due next Friday 9/24, at 8pm.
- Class pacing: Just right – 53/87, A bit too fast – 29/87.
- If things feel way too fast, come to my office hours or send me an email and we can discuss how to make things more manageable.
- **Reminder:** My office hours are now Thursdays 9am on Zoom. Pratheba's in person office hours are Mondays 12pm-1pm in CS 207.

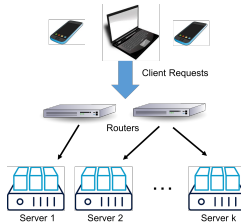
Last Class:

- 2-Level Hashing Analysis (linearity of expectation and Markov's inequality)
- 2-universal and pairwise independent hash functions
- Start on hashing for load balancing.

This Time:

- Finish hashing for load balancing. 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.

Randomized Load Balancing:



- n requests randomly assigned to k servers.
- Expected load on server i is $\mathbb{E}[R_i] = \frac{n}{k}$.
- By Markov's inequality, if we provision each server to handle twice this expected load (so $\frac{2n}{k}$ requests), it will be overloaded with probability $\leq 1/2$.

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:

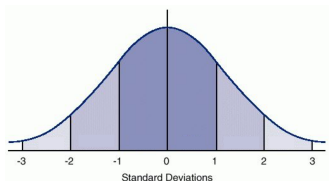
$$\Pr(|X - \mathbb{E}[X]| \geq t) = \Pr(X^2 \geq t^2) \leq \frac{\mathbb{E}[X^2]}{t^2} \frac{\text{Var}[X]}{t^2}.$$

(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?



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

Why is this so powerful?

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

Consider drawing independent identically distributed (i.i.d.) random variables $\mathbf{X}_1, \dots, \mathbf{X}_n$ with mean μ and variance σ^2 .

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

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

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

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

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.

$$\begin{aligned} \text{Var}[R_{i,j}] &= \mathbb{E} \left[(R_{i,j} - \mathbb{E}[R_{i,j}])^2 \right] \\ &= \Pr(R_{i,j} = 1) \cdot (1 - \mathbb{E}[R_{i,j}])^2 + \Pr(R_{i,j} = 0) \cdot (0 - \mathbb{E}[R_{i,j}])^2 \\ &= \frac{1}{k} \cdot \left(1 - \frac{1}{k}\right)^2 + \left(1 - \frac{1}{k}\right) \cdot \left(0 - \frac{1}{k}\right)^2 \\ &= \frac{1}{k} - \frac{1}{k^2} \leq \frac{1}{k} \implies \text{Var}[R_i] \leq \frac{n}{k}. \end{aligned}$$

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 CHEBYSHEVS

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 extremely 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?

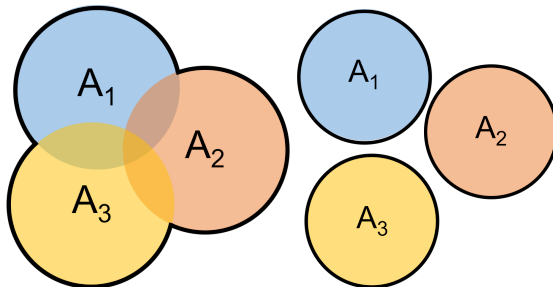
$$\Pr \left(\max_i (\mathbf{R}_i) \geq \frac{2n}{k} \right) = \Pr \left(\left[\mathbf{R}_1 \geq \frac{2n}{k} \right] \cup \left[\mathbf{R}_2 \geq \frac{2n}{k} \right] \cup \dots \cup \left[\mathbf{R}_k \geq \frac{2n}{k} \right] \right) = \Pr$$

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 . $\mathbb{E}[\mathbf{R}_i] = \frac{n}{k}$. $\text{Var}[\mathbf{R}_i] = \frac{n}{k}$.

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.

On the first problem set, you will prove the union bound, as a consequence of Markov's inequality.

APPLYING THE UNION BOUND

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?

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

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

ANOTHER VIEW ON THIS PROBLEM

The number of servers must be small compared to the number of requests ($k = O(\sqrt{n})$) for the maximum load to be bounded in comparison to the expected load with good probability.

- There are many requests routed to a relatively small number of servers so the load seen on each server is close to what is expected via law of large numbers.
- **A More Natural Variant:** Given n requests, and assuming all servers have fixed capacity C , how many servers should you provision so that with probability $\geq 99/100$ no server is assigned more than C requests?

n : total number of requests, k : number of servers randomly assigned requests.

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

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

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

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

\mathbf{H} has a simple Binomial distribution, so can compute these probabilities exactly.

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