COMPSCI ~~514~~ 614 ~~exam~~: Randomized Algorithms ~~and Probabilistic~~ Data Analysis

w/ Apps to

Prof. Cameron Musco

University of Massachusetts Amherst. Spring 2024

Lecture 4

## Logistics

- Problem Set 1 is due next Wednesday 2/21 at 11:59pm.
- Most people think the lectures are 'just right' or 'a bit too fast'. I'll try to slow down a bit. If you feel that you are really falling behind, let me know.
- If you are confused on something please ask about it – certainly you are not the only one!

## Summary

$$Pr(X \geq t) \leq \frac{E[X]}{t}$$

$$Pr(|X - EX| \geq t) \leq \frac{var(X)}{t^2}$$

Last Time:

- Concentration bounds – Markov's and Chebyshev's inequalities.
- The union bound. $Pr(A_1 \cup \ldots \cup A_n) \leq \sum P(A_i)$
- Coupon collecting, statistical estimation.
- Randomized load balancing and ball-into-bins

Last Time:

- Concentration bounds – Markov's and Chebyshev's inequalities.

- The union bound.

- Coupon collecting, statistical estimation.

- Randomized load balancing and ball-into-bins

Today:

- Stronger concentration bounds for sums of independent random variables. I.e., exponential concentration bounds.

- Applications to balls-into-bins and linear probing analysis.

*Central limit theorem*

error $\approx \frac{1}{\sqrt{n}}$

**Question 4**

Not complete

Points out of 1.00

⚑ Flag question

✏ Edit question

Let's say I have two biased coins -- one hits heads with probability $1/2 + \epsilon$ and tails with probability $1/2 - \epsilon$. The other hits tails with probability $1/2 + \epsilon$ and heads with probability $1/2 - \epsilon$.  \multiarm bandit , hypothesis testing

How many independent flips of the coins must I perform to distinguish them from each other with probability at least $2/3$.

- a. $O(\log(1/\epsilon))$
- b. $O(1/\epsilon)$
- c. $O(1/\epsilon^2)$
- d. $O(1/\epsilon^4)$

[Check]

$C_H = \#$ times heads hits heads

$C_T = \#$ tails coin

flip tiies

$\mathbb{E}[C_H - C_T] = 2\epsilon n \geq 0$

$\mathbb{E}[C_H] - \mathbb{E}[C_T] = \frac{n}{2} + \epsilon n - \left(\frac{n}{2} - \epsilon n\right)$

$Pr\left(C_H > C_T\right)$

$Pr\left(C_H - C_T > 0\right)$

$Var[C_H - C_T] = Var[C_H] + Var[C_T]$

$Pr(C_H - C_T < 0) \leq$

$Pr\left([(C_H - C_T) - \mathbb{E}(C_H - C_T)] \geq 2\epsilon n\right) \leq \frac{n/2}{4\epsilon^2 n^2} = \frac{1}{8\epsilon^2 n} \leq \frac{2}{9}$

$Var(C_H - C_T) \leq \frac{n}{2}$

$Var\left[\sum_{i=1}^{n} C_{H,i}\right]$

$= n \cdot Var(C_{H,i})$

$\leq \frac{1}{4}$

$p(1-p) = \left(\frac{1}{2} + \epsilon\right)\left(\frac{1}{2} - \epsilon\right)$

$= \frac{1}{4} - \epsilon^2$

$$Var(X+Y) \leq 2Var(X) + 2Var(Y)$$

$X = 0 \ \text{w.p } 1-p$    either $n$ or $-n$    $n.p \ 2p$    $Var(x) = E[(X-Ex)^2]$

$= (1-p) \cdot 0$

**Question 5**

Not complete

Points out of 1.00

⚐ Flag question

✱ Edit question

You roll a fair 6-sided die $n$ times independently. You look at the difference between the number of times you rolled a "1" the number of times you rolled a "2". Roughly, how big do we expect this difference to be in magnitude? **Hint:** What is the variance of this difference?

- a. $\Theta(n)$
- b. $\Theta(\sqrt{n})$
- c. $\Theta(\log n)$
- d. $\Theta\left(\frac{\log n}{\log \log n}\right)$

Check

$+ 2p \cdot n^2$

$Var(x) \leq p \cdot n^2$

s.d: $\sqrt{p} \cdot n$    $E[|x|] = 2p \cdot n$

$X_1 = \# \ ones$

$X_2 = \# \ twos$

$Var(X_1 - X_2)$    $\leq \sqrt{n}$

$X_1 - X_2$

$X_1 - X_2 = \sum_{i=1}^{n} D_i$     $D_i = 1 \ \text{w.p} \ 1/6$

$= -1 \ \text{w.p} \ 1/6$

$Var(X_1 - X_2) \leq Var(X_1) + Var(X_2) \ ?$     $0 \ \text{w.p} \ 2/3$

$Var(X_1 - X_2) = Var(\sum D_i) \leq n = E[(X_1 - X_2)^2]$

$E[|X_1 - X_2|]$     $var(D_i) \leq 1$

# Balls Into Bins

# Balls Into Bins

I throw $m$ balls independently and uniformly at random into $n$ bins. What is the maximum number of balls any bin?



Bin 1   Bin 2   Bin 3

- Applications to randomized load balancing
- Analysis of hash tables using chaining.

I throw $m$ balls independently and uniformly at random into $n$ bins. What is the maximum number of balls any bin?

$b_1 = 3$

$b_2 = 2$

$b_3 = 2$

Bin 1

Bin 2

Bin 3

- Applications to randomized load balancing
- Analysis of hash tables using chaining.
- **Direct Proof:** For any bin $i$, $\Pr[b_i \geq \frac{c \ln n}{\ln \ln n}] \leq \frac{1}{n^{c-o(1)}}$. Thus, via union bound, the maximum load is exceeds $\frac{c \ln n}{\ln \ln n}$ with probability at most $\frac{1}{n^{c-1-o(1)}}$.

$$n \cdot \frac{1}{n^{c-o(1)}}$$

# Balls Into Bins Via Chebyshev's Inequality

In our balls into bins analysis we directly bound
$\Pr[b_i \geq k] \leq \left(\frac{e}{k}\right)^k \cdot \frac{1}{1-e/k}$.

$n = m$

**Think Pair Share:** Give an upper bound on this probability using Chebyshev's inequality. Hint: write $b_i$ as a sum of $n$ indicator random variables and compute $\text{Var}[b_i]$ and/or $\mathbb{E}[b_i^2]$.

$$\Pr(b_i \geq k) \leq O\left(\frac{1}{k^2}\right)$$

$$b_i = \sum_{j=1}^{n} X_j \qquad \text{Var}(b_i) = \sum \text{Var}(X_j) = n \cdot \text{Var}(X_j)$$

$$\text{Var}(b_i) \leq 1 \qquad \overset{\leq 1}{\text{Var}(b_i)} = \mathbb{E} b_i^2 - [\mathbb{E} b_i]^2 \qquad \frac{1}{n}\left(1 - \frac{1}{n}\right) \leq \frac{1}{n}$$

$$\mathbb{E}(b_i) = 1 \qquad \mathbb{E} b_i^2 \leq 1+1 \leq 2$$

$$\Pr[b_i \geq k] \leq \Pr[|b_i - \mathbb{E} b_i| \geq k-1] \leq \frac{\text{Var}(b_i)}{(k-1)^2} \leq \frac{1}{(k-1)^2}$$

$$\Pr[b_i \geq k] = \Pr(b_i^2 \geq k^2) \leq \frac{2}{k^2}$$

## Balls Into Bins Via Chebyshev's Inequality

By **Chebyshev's Inequality:** $\Pr\left[b_i \geq k\right] \leq \frac{2}{k^2}$.

Setting $k = c\sqrt{n}$, $\Pr\left[b_i \geq c\sqrt{n}\right] \leq \frac{2}{c^2 n}$. So via a union bound:

$$\Pr\left[\underbrace{\max_{i=1,\ldots,n} b_i \geq c\sqrt{n}}\right] \leq n \cdot \frac{2}{c^2 n} \leq \frac{2}{c^2}.$$

## Balls Into Bins Via Chebyshev's Inequality

By Chebyshev's Inequality: $\Pr\left[\mathbf{b}_i \geq k\right] \leq \frac{2}{k^2}$.

Setting $k = c\sqrt{n}$, $\Pr\left[\mathbf{b}_i \geq c\sqrt{n}\right] \leq \frac{2}{c^2 n}$. So via a union bound:

$$\Pr\left[\max_{i=1,\ldots,n} \mathbf{b}_i \geq c\sqrt{n}\right] \leq n \cdot \frac{2}{c^2 n} \leq \frac{2}{c^2}.$$

**Upshot:** Chebyshev's inequality bounds the maximum load by $O(\sqrt{n})$ with good probability, as compared to $O\left(\frac{\log n}{\log \log n}\right)$ for the direct proof. It is quite loose here.

$\mathrm{Var}(X+Y) = \mathrm{Var}(X) + \mathrm{Var}(Y) + 2\,\mathbb{E}(X - \mathbb{E}X)(Y - \mathbb{E}Y)$

$\leq 2\sqrt{\mathrm{Var}(X)\mathrm{Var}(Y)}$  AM-GM

$Y = X$

$\mathrm{Var}(X+Y) \leq 2\mathrm{Var}(X) + 2\mathrm{Var}(Y)$   $\hookleftarrow \leq \mathrm{Var}(X) + \mathrm{Var}(Y)$

$\mathrm{Var}(2X) = 4\mathrm{Var}(X)$   Can't get lower bound b/c we could have $X = -Y$

8

By Chebyshev's Inequality: $\Pr\left[\mathbf{b}_i \geq k\right] \leq \frac{2}{k^2}$.

Setting $k = c\sqrt{n}$, $\Pr\left[\mathbf{b}_i \geq c\sqrt{n}\right] \leq \frac{2}{c^2 n}$. So via a union bound:

$$\Pr\left[\max_{i=1,\dots,n} \mathbf{b}_i \geq c\sqrt{n}\right] \leq n \cdot \frac{2}{c^2 n} \leq \frac{2}{c^2}.$$

**Upshot:** Chebyshev's inequality bounds the maximum load by $O(\sqrt{n})$ with good probability, as compared to $O\left(\frac{\log n}{\log \log n}\right)$ for the direct proof. It is quite loose here.

Chebyshev's and Markov's inequalities are extremely valuable because they are very general – require few assumptions on the underlying random variable. But by using assumptions, we can often get tighter analysis.

Exponential Concentration Bounds

Markov's Inequality: $\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$. First moment.

$t = \sqrt{n} c$

Chebyshev's Inequality: $\Pr[X \geq t] \leq \frac{\mathbb{E}[X^2]}{t^2}$. Second moment.

$$Pr(X \geq t) = Pr(X^2 \geq t^2) \leq \frac{\mathbb{E}[X^2]}{t^2} = \frac{n}{t^2} = \frac{1}{c^2}$$

$$= Pr(X^4 \geq t^4) \leq \frac{\mathbb{E}[X^4]}{t^4} \leq \frac{n^2}{t^4} = \frac{1}{c^4}$$

$$X = \sum_{i=1}^{n} X_i$$

$X_i = 1 \text{ w.p } \frac{1}{2}$

$X_i = -1 \text{ w.p. } 1/2$

$$\frac{\mathbb{E}X^6}{t^6} \leq \frac{n^3}{t^6}$$

$\mathbb{E}X = 0$

$\mathbb{E}X^2 = Var(X) = \sum Var(X_i) = n \cdot 1$

$$\mathbb{E}X^4 = \mathbb{E}\left[\left(\sum_{i=1}^{n} X_i\right)^4\right] = n \cdot \underbrace{\sum \mathbb{E}X_i^4}_{n} + \binom{4}{2} \underbrace{\sum \mathbb{E}X_i^2 \mathbb{E}X_j^2}_{n^2} + \text{other stuff}$$

$= \mathbb{E}[(X_1 + \cdots X_n)^4]$

$\underset{=}{n^2}$

$X_i^2 X_j X_k$

$X_j^3 X_i$

0

9

## Higher Moments

**Markov's Inequality:** $\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$. First moment.

**Chebyshev's Inequality:** $\Pr[X \geq t] \leq \frac{\mathbb{E}[X^2]}{t^2}$. Second moment.

Often (not always!) we can obtain tighter bounds by looking to higher moments of the random variable.

## Higher Moments

**Markov's Inequality:** $\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$. First moment.

**Chebyshev's Inequality:** $\Pr[X \geq t] \leq \frac{\mathbb{E}[X^2]}{t^2}$. Second moment.

Often (not always!) we can obtain tighter bounds by looking to higher moments of the random variable.

**Moment Generating Function:** Consider for any $z > 0$:

$$M_z(X) = e^{z \cdot X} = \sum_{k=0}^{\infty} \frac{z^k X^k}{k!}$$

by

## Higher Moments

Markov's Inequality: $\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$. First moment.

Chebyshev's Inequality: $\Pr[X \geq t] \leq \frac{\mathbb{E}[X^2]}{t^2}$. Second moment.

Often (not always!) we can obtain tighter bounds by looking to higher moments of the random variable.

Moment Generating Function: Consider for any $z > 0$:

$$M_z(X) = e^{z \cdot X} = \sum_{k=0}^{\infty} \frac{z^k X^k}{k!}$$

$e^{z \cdot t}$ is non-negative, and monotonic for any $z > 0$. So can bound via Markov's inequality, $\Pr[X \geq t] = \Pr[M_z(X) \geq e^{zt}] \leq \frac{\mathbb{E}[M_z(X)]}{e^{zt}}$.

## Higher Moments

**Markov's Inequality:** $\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$. First moment.

**Chebyshev's Inequality:** $\Pr[X \geq t] \leq \frac{\mathbb{E}[X^2]}{t^2}$. Second moment.

Often (not always!) we can obtain tighter bounds by looking to higher moments of the random variable.

**Moment Generating Function:** Consider for any $z > 0$:

$$M_z(X) = e^{z \cdot X} = \sum_{k=0}^{\infty} \frac{z^k X^k}{k!}$$

$e^{z \cdot t}$ is non-negative, and monotonic for any $z > 0$. So can bound via Markov's inequality, $\Pr[X \geq t] = \Pr[M_z(X) \geq e^{zt}] \leq \frac{\mathbb{E}[M_z(X)]}{e^{zt}}$.

By appropriately picking $z$ and bounding $\mathbb{E}[M_z(X)]$, we can obtain a variety of exponential tail bounds. Typically require that X is a sum of bounded and independent random variables

$\vdash$ Hoeffding, Chernoff, Berstein, Azuma H...
Berry Esseen.

## The Chernoff Bound

**Chernoff Bound (simplified version):** Consider <u>independent</u> random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$ and let $X = \sum_{i=1}^{n} X_i$. Let $\mu = \mathbb{E}[X] = \mathbb{E}[\sum_{i=1}^{n} X_i]$. For any $\delta \geq 0$

$$\Pr(X \geq (1+\delta)\mu) \leq \frac{e^{\delta\mu}}{(1+\delta)^{(1+\delta)\mu}}$$

$\delta \geq 5$

not neccesarily binomial

$$\leq \frac{e^{5\mu}}{6^{6\mu}} \qquad \delta = 5$$

$$\leq \frac{6^{5\mu}}{6^{6\mu}} \leq \frac{1}{6^{\mu}}$$

## The Chernoff Bound

**Chernoff Bound (simplified version):** Consider independent random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$ and let $X = \sum_{i=1}^{n} X_i$. Let $\mu = \mathbb{E}[X] = \mathbb{E}[\sum_{i=1}^{n} X_i]$. For any $\delta \geq 0$

$$\Pr\left(X \geq (1+\delta)\mu\right) \leq \frac{e^{\delta\mu}}{(1+\delta)^{(1+\delta)\mu}}$$

**Chernoff Bound (alternate version):** Consider independent random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$ and let $X = \sum_{i=1}^{n} X_i$. Let $\mu = \mathbb{E}[X] = \mathbb{E}[\sum_{i=1}^{n} X_i]$. For any $\delta \geq 0$

$$\Pr\left(\left|\sum_{i=1}^{n} X_i - \mu\right| \geq \delta\mu\right) \leq 2\exp\left(-\frac{\delta^2\mu}{2+\delta}\right).$$

As $\delta$ gets larger and larger, the bound falls off exponentially fast.

## Balls Into Bins Via Chernoff Bound

Recall that $b_i$ is the number of balls landing in bin $i$, when we randomly throw $n$ balls into $n$ bins.

- $b_i = \sum_{i=1}^{n} I_{i,j}$ where $I_{i,j} = 1$ with probability $1/n$ and 0 otherwise. $I_{i,1}, \ldots I_{i,n}$ are independent.

### Balls Into Bins Via Chernoff Bound

Recall that $b_i$ is the number of balls landing in bin $i$, when we randomly throw $n$ balls into $n$ bins.

- $b_i = \sum_{i=1}^{n} I_{i,j}$ where $I_{i,j} = 1$ with probability $1/n$ and 0 otherwise. $I_{i,1}, \ldots I_{i,n}$ are independent.
- Apply Chernoff bound with $\mu = \mathbb{E}[b_i] = 1$:

$$\Pr[b_i \geq k] \leq \frac{e^k}{(1+k)^{(1+k)}}.$$

## Balls Into Bins Via Chernoff Bound

Recall that $\mathbf{b}_i$ is the number of balls landing in bin $i$, when we randomly throw $n$ balls into $n$ bins.

- $\mathbf{b}_i = \sum_{i=1}^{n} \mathbf{I}_{i,j}$ where $\mathbf{I}_{i,j} = 1$ with probability $1/n$ and 0 otherwise. $\mathbf{I}_{i,1}, \ldots \mathbf{I}_{i,n}$ are independent.
- Apply Chernoff bound with $\mu = \mathbb{E}[\mathbf{b}_i] = 1$:

$$\Pr[\mathbf{b}_i \geq k] \leq \frac{e^k}{(1+k)^{(1+k)}}.$$

- For $k \geq \frac{c \log n}{\log \log n}$ we have:

$$\Pr[\mathbf{b}_i \geq k] \leq \frac{e^{\frac{c \log n}{\log \log n}}}{\left( \frac{c \log n}{\log \log n} \right)^{\frac{c \log n}{\log \log n}}} =$$