

COMPSCI 614: Randomized Algorithms with Applications to Data Science

Prof. Cameron Musco

University of Massachusetts Amherst. Spring 2024.

Lecture 13 (Midterm Review)

Midterm Logistics

- Thursday in class 10-11:15 am.
- Closed book, no calculator.
- Study guide posted under course schedule on the midterm line.
- Practice midterm posted in Moodle.

Midterm Format/Content

Rough Format:

- 1. Set of true/false or always/sometimes/never questions
- 2. Set of short answer questions.
- 3-4. Two multi-part 'problem set like' questions.
- 5. One bonus question.

Midterm Format/Content

Rough Format:

- 1. Set of true/false or always/sometimes/never questions
- 2. Set of short answer questions.
- 3-4. Two multi-part 'problem set like' questions.
- 5. One bonus question.

Content:

- See study guide for a detailed list.
- Should know key tools/analysis approaches. But don't need to memorize derivations from class.
- It will help to be familiar with homework problems, but don't need to memorize answers/derivations from them.

Midterm Studying

- I would focus on doing practice problems more than reviewing material.
- Definitely do the practice exam in semi-timed, closed note environment.
- Try to do as many other practice questions from the study guide, by redoing homework questions/problems given in class, etc.

Practice

A parallel computer consists of n processors and n memory modules. During a step, each processor sends a memory request to a random module. A module that receives 1 or 2 requests satisfies its requests; modules that receive more than two requests will not satisfy them. What is the expected number of modules that satisfy their requests?

Practice

A parallel computer consists of n processors and n memory modules. During a step, each processor sends a memory request to a random module. A module that receives 1 or 2 requests satisfies its requests; modules that receive more than two requests will not satisfy them. What is the expected number of modules that satisfy their requests?

$X = \#$ modules satisfying requests

$$\mathbb{E}[X] = \sum_{i=1}^n \mathbb{E}[X_i]$$

$$X = X_1 + \dots + X_n$$

$X_i = 1$ if module i satisfies
 0 otherwise

$b_i = \#$ balls in bin i

$$\Pr(b_i=1) + \Pr(b_i=2)$$

$$= n \cdot \frac{1}{n} \left(1 - \frac{1}{n}\right)^{n-1} + \frac{\binom{n}{2}}{n^2} \left(1 - \frac{1}{n}\right)^{n-2} \rightarrow \mathbb{E}[X] = n \left(1 - \frac{1}{n}\right)^{n-1} + \frac{\binom{n}{2}}{n} \left(1 - \frac{1}{n}\right)^{n-2}$$

Practice

You store n items in a hash table with $2n$ buckets using a fully random hash function. Give an upper bound on the maximum load in any bucket, which holds with probability at least $1 - 1/n$. Does the answer change significantly if you have n instead of $2n$ buckets?

$$E[X] = n \left(1 - \frac{1}{2n}\right)^{n-1} + \frac{\binom{n}{2}}{n} \left(1 - \frac{1}{2n}\right)^{n-2}$$

$$n \rightarrow \infty \quad \approx \frac{n}{e} + \frac{n(n-1)}{2n} \cdot \frac{1}{e} = n \left(\frac{1}{e} + \frac{1}{2e} \right)$$

Practice

$b_1 \dots b_n$

You store n items in a hash table with $2n$ buckets using a fully random hash function. Give an upper bound on the maximum load in any bucket, which holds with probability at least $1 - 1/n$. Does the answer change significantly if you have n instead of $2n$ buckets?

$$\mathbb{E} b_i = \frac{1}{2}$$

$$\begin{aligned} \max \text{ load} &\approx O\left(\frac{\log n}{\log \log n}\right) \\ \max \text{ load} &\leq O(\log n) \end{aligned}$$

$O(\sqrt{n})$
 $O(n)$

$$\begin{aligned} \Pr(\max \text{ load} \geq t) &= \Pr(b_1 \geq t \text{ or } b_2 \geq t \dots \text{ or } b_n \geq t) \\ &\stackrel{\text{union bound}}{\leq} n \cdot \Pr(b_i \geq t) \end{aligned}$$

we need to set t s.t.

$$\Pr(b_i \geq t) \leq \frac{1}{n^2}$$

$$\begin{aligned} \text{Then } \Pr(\max \text{ load} \leq t) &\geq 1 - n \Pr(b_i \geq t) \\ &\geq 1 - \frac{n}{n^2} \geq 1 - \frac{1}{n} \end{aligned}$$

Practice

Consider the scenario above, where you use linear probing instead of chaining. Does the expected look up time change significantly if you use n rather than $2n$ buckets?

$$\Pr(b \geq t)$$

$$\Pr\left(\sum_{j=1}^n b_j \geq t\right)$$

where $b_j = 1$ if bucket j lands in bin
 $b_j = 0$ otherwise

$$\Pr\left(\sum b_j \geq \delta n\right) \leq \exp\left(\frac{-\mu \delta^2}{2t\delta}\right)$$

$$\mu = \mathbb{E}b_j \quad \mu = \frac{1}{2} \quad \delta = 2t$$

$$\Pr\left(\sum b_j \geq t\right) \leq \exp\left(\frac{-\frac{1}{2} \cdot 2t^2}{2t \cdot 2t}\right) \approx \exp(-t)$$

$t = \frac{\log n}{\exp(-t)} \leq \frac{1}{n^2}$ for constant c

Practice

You have an algorithm that succeeds with probability $2/3$. You run it t times independently and would like to ensure that the probability that the algorithm fails on a majority (i.e., $> t/2$) of these trials is at most δ . How large should you set t ?

Practice

You have an algorithm that succeeds with probability $2/3$. You run it t times independently and would like to ensure that the probability that the algorithm fails on a majority (i.e., $> t/2$) of these trials is at most δ . How large should you set t ?

Practice

$$Y = \text{# failures} \quad E[Y] = t/3$$

You have an algorithm that succeeds with probability $2/3$. You run it t times independently and would like to ensure that the probability that the algorithm fails on a majority (i.e., $> t/2$) of these trials is at most δ . How large should you set t ?

$$\Pr(X - EX \geq \frac{t}{6}) \leq \frac{\text{Var}(X)}{\frac{t^2}{36}} \leq \frac{\frac{t}{3}}{\frac{t^2}{36}} = \frac{12}{t}$$

→ median

→ BPP success prob can be boosted

→ Chernoff $t = O(\log(1/\delta))$

→ Bernstein

$$t = O(1/\delta)$$

→ Chebyshev

→ Markov

$$\Pr(Y \geq \frac{t}{2}) = \Pr(Y \geq E(Y) \cdot \frac{3}{2}) \leq \frac{2}{3}$$

Practice

How to improve dependence on d : use Bernstein
 - use median trick $O(\frac{\log(1/\delta)}{\epsilon^2})$

There is a slot machine that on each pull pays out a random prize between $[0, 1]$. After t plays roughly how good of an estimate of the expected value of playing do you have?

$$\mathbb{E}[X] = 0$$

$$\mathbb{E}[X] = 1/2$$

$$1/\sqrt{T}$$

$$\text{Var}[X] \leq 1$$

to get ϵ error in estimating value w.p. $1-\delta$ how big should T be?

$$\bar{X} = \frac{1}{T} \sum_{i=1}^T X_i$$

$$\text{Var}(\bar{X}) \leq \frac{1}{T} \sum_{i=1}^T \text{Var}(X_i)$$

$$T = O\left(\frac{1}{\delta \epsilon^2}\right)$$

$$\leq \frac{1}{T}$$

$$\text{Var}(\bar{X}) = \mathbb{E}(\bar{X} - \mathbb{E}\bar{X})^2 = \frac{1}{T}$$

by Chebyshev's:

$$\Pr(|\bar{X} - \mathbb{E}\bar{X}| \geq \epsilon) \leq \frac{\text{Var}(\bar{X})}{\epsilon^2} \leq \frac{1}{T \epsilon^2}$$

Practice

There is a slot machine that on each pull pays out a random prize between $[0, 1]$. After t plays roughly how good of an estimate of the expected value of playing do you have?

Practice

$$\text{Var}(X) = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$$

You would like to estimate the inner product of a vector $a \in \mathbb{R}^n$ with a vector $x \in \mathbb{R}^n$. You know a completely but don't know x . Describe how to obtain an estimate of the inner product without reading all of x via important sampling.

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n a_i x_i$$

$$M = \sum a_i x_i$$

sample ij w.p. P_{ij}
how to choose P_{ij} ?

$$M = \frac{1}{\dagger} \sum_{j=1}^{\dagger} a_{ij} x_{ij} \cdot \frac{1}{P_{ij}}$$

$$\text{Var}(M) = \frac{1}{\dagger} \text{Var}(a_{ij} x_{ij} \cdot \frac{1}{P_{ij}}) = \frac{1}{\dagger} \mathbb{E} (a_{ij} x_{ij} \frac{1}{P_{ij}} - \mathbb{E}[\dots])^2 \leq \frac{1}{\dagger} \mathbb{E} (a_{ij} x_{ij} \frac{1}{P_{ij}^2})$$

Practice

You would like to estimate the inner product of a vector $a \in \mathbb{R}^n$ with a vector $x \in \mathbb{R}^n$. You know a completely but don't know x . Describe how to obtain an estimate of the inner product without reading all of x via important sampling.

$$\text{Var}(m) \leq \frac{1}{T} \mathbb{E} \left[a_{i_j}^2 x_{i_j}^2 \frac{1}{p_{i_j}^2} \right]$$

$$\downarrow$$
$$= \sum_{j=1}^n p_j \cdot a_j^2 x_j^2 \frac{1}{p_j^2}$$

$$\text{Var}(m)$$

$$= \frac{1}{T} \sum_{j=1}^n \frac{a_j^2 x_j^2}{p_j}$$

$$= \frac{1}{T} \sum_{j=1}^n \frac{a_j^2 x_j^2}{\frac{a_j^2 x_j^2}{\|a\|_2^2}} = \frac{1}{T} \|a\|_2^2 \|x\|_2^2$$

$$\text{opt: } p_j \propto \frac{a_j x_j}{\sum a_j x_j}$$
$$p_j = \frac{a_j^2}{\sum a_j^2} = \frac{a_j^2}{\|a\|_2^2}$$

$$\|a\|_\infty, \|x\|_\infty \leq M$$

$$\min V = \frac{1}{T} \sum_{j=1}^n \frac{a_j^2 x_j^2}{p_j} = \frac{n}{T} \sum a_j^2 x_j^2 = \frac{n^2 \cdot m^4}{T}$$

$\left(\frac{n \cdot m^2}{\sqrt{T}} \right)^2$

$$\frac{\partial V}{\partial p_1} = \dots = \frac{\partial V}{\partial p_n}$$

$$\text{s.t. } p_j \in [0, 1]$$

$$\underline{\underline{\sum p_j = 1}}$$

$$\left(\approx \frac{a_j^2 x_j^2}{p_j^2} \right)$$

$$\text{so set } p_j^2 \propto a_j^2 x_j^2$$

what could you get w/ uniform sampling?

$$x = \begin{bmatrix} \infty \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} \infty \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbb{E} [\underbrace{\|AB - \bar{C}\|_F^2}]$$

