# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.
Lecture 13 (Midterm Review)

## Summary

Last Class:

- Introduced the idea of low-distortion embeddings and the JL Lemma.

$x - y$    $\|x - y\|_2$          $\|y\|$

- Reduction of JL Lemma to the Distributional JL Lemma via union bound.

- We will finish the proof of the JL Lemma after the midterm. Ignore any practice questions on this topic.

This Class:

- Midterm review.

**Rough Outline:** (subject to small changes)

*1 pt for answer*

*1 pt for explanation*

25/30

- Question 1: 4-5 always, sometimes, nevers or (true falses.)
- Question 2: 3-4 short answers, sort of like quiz questions.
- Question 3-4: Multipart questions, similar to core competency problems.

6

- Question 5: Extra credit question. Similar to a harder core competency problem.

3)/30

3

# Questions

Content, Format, or Logistics Questions?

- Bernstein
- Chernoff
- bloom filter FPR

$$\binom{m}{2} \leq \frac{m^2}{2}$$

$$m^2$$

midterm concepts.pdf

$$Var(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2$$
$$= p - p^2$$
$$\{0,1\} \quad 0: 1-p$$
$$1: p$$

$$(1+x) < e^x$$

$$\frac{\exp(\bar{r}_j, k)}{} \leq \frac{1}{k}$$

# Questions

$$\sum_{i=1}^{n} \underline{\mathbb{E} S_i^2}$$

$X_1$

$$S_1 = 2 \quad \circ \; \circ \; \circ$$

$$S_2 = 3 \quad \circ \; \circ \; \circ$$

$$S_3 = 0$$

$$\circ$$

$$\circ$$

$\vdots$

$\vdots$

$X_n$

$$S_n \quad \circ \; \circ \; \circ \; \circ r$$

$$\left[ \; \vdots \; \left\{ \begin{array}{c} S_1^2 = 4 \\ \\ \end{array} \right. \right.$$

$$\left\{ \begin{array}{c} \vdots \\ \circ \end{array} \right\} S_2^2$$

$$\begin{array}{c} \circ \\ \circ \\ \circ \\ \vdots \\ \circ \end{array} \quad S_n^2$$

no collision

$$S_i^2 = \left( \sum_j X_{ij} \right)^2$$

1 if item
$j$ hashes
to bucket $i$

$$\mathbb{E} \sum_k \sum_j X_{ij} \cdot X_{ik}$$

$$\sum_k \sum_j \underline{\mathbb{E} X_{ij} X_{ik}}$$

pairwise ind

$$\frac{S_i^2}{2S_i^2}$$

$$\left( \sum_k \mathbb{E} X_{ik} \quad \sum_{i \neq j} \sum \mathbb{E} X_{ik} X_{jk} \right)$$

$$\sum_i \sum_k \mathbb{E} X_{ik}$$

2-universal

$$Pr(h(x_j) = h(y_k)) \cdot n$$

$$n$$

$$\sum \mathbb{E} S_i^2$$

$$\frac{1}{n} \sum_{i=1}^{n}$$

collisions

$$\frac{1}{n^2} = Pr\left( \begin{array}{c} h(x_j) = i \text{ AND} \\ h(x_k) = i \end{array} \right)$$

5

Random Hash Functions

$h(x): U_2 \to [n]$

$Pr(h(x) = h(y)) \leq \frac{1}{n}$

$Pr(h(x) = h(y) = 1) = 0$

$Pr(h(x) = h(y) = 2) = \frac{2}{n}$

LSH

Fully random
hash function

pairwise

$h(x) = rand(n)$

$h(x), h(y), h(z), \dots$

are independent

2-universal $Pr(h(x) = h(y)) \leq \frac{1}{n}$

pairwise independent

$Pr(h(x) = i \text{ AND } h(y) = j) = \frac{1}{n^2}$

$h(2) = h(1) + 1 \mod (n)$

LSH

$Pr(h(1) = 1 \ and \ h(2) = 1)) = 0$

hash tables

rarely

CMS

$Pr(h(x) = h(y))$

$= \sum_{i=1}^{n} Pr(h(x) = i \text{ and } h(y) = i)$

$x \neq y$

$\sum_{i=1}^{n} \frac{1}{n^2} = \frac{1}{n}$

6

# Concentration Bounds



**Concentration Bound Requirements**

$Var(X_i) = p - p^2$
$p(1-p)$
$\leq \frac{1}{4}$

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

| Markov's | Chebyshev's | Chernoff | Bernstein |
|---|---|---|---|
| $E[X]$ | $Var(X)$ | $X = \sum X_i$ in 2. | $X = \sum X_i$ in 2. |
| $X > 0$ | $E(X)$ | $X_i \in \{0,1\}$ | $|X_i| \leq M$ |
| | | $E[X] = \mu$ | $\sigma^2 = Var(X)$ |
| | | $Pr(|X - E X| > \delta M)$ | $t$ |

$m = 1 \quad m = 2$

$$Pr(|X - E X| > t) \leq \frac{Var(X)}{t^2}$$

$\frac{1}{10}$

$$exp\left(\frac{-\frac{1}{2} \cdot \mu}{2t + \frac{1}{2}}\right)$$

$3$

$\delta = \frac{10}{\mu}$

7

3. Consider an algorithm $\mathcal{A}$ running in time $T(\mathcal{A})$, that with probability .6 outputs an estimate of the number of triangles in an input graph up to error $\pm 100$, and with probability .4 outputs some bad estimate with worse error. Describe an algorithm that outputs an estimate of the number of triangles in an input graph up to error $\pm 100$ with probability $\geq .99$ and runs in time $O(T(\mathcal{A}))$.

The Chernoff bound states that for independent random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$, letting $\mu = \mathbb{E}\left[\sum_{i=1}^{n} X_i\right]$, for any $\delta > 0$,

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

3. Consider an algorithm $\mathcal{A}$ running in time $T(\mathcal{A})$, that with probability .6 outputs an estimate of the number of triangles in an input graph up to error $\pm 100$, and with probability .4 outputs some bad estimate with worse error. Describe an algorithm that outputs an estimate of the number of triangles in an input graph up to error $\pm 100$ with probability $\geq .99$ and runs in time $O(T(\mathcal{A}))$.

The Chernoff bound states that for independent random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$, letting $\mu = \mathbb{E}\left[\sum_{i=1}^{n} X_i\right]$, for any $\delta > 0$,

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

# Example Problems

2. Assume there are 1000 registered users on your site $u_1, \ldots, u_{1000}$, and in a given day, each user visits the site with some probability $p_i$. The event that any user visits the site is independent of what the other users do. Assume that $\sum_{i=1}^{1000} p_i = 500$.

    (a) Let $\mathbf{X}$ be the number of users that visit the site on the given day. What is $\mathbb{E}[\mathbf{X}]$.

    (b) Apply a Chernoff bound to show that $\Pr[\mathbf{X} \geq 600] \leq .01$.

    (c) Apply Markov's inequality and Chebyshev's inequality to bound the same probability. How do they compare?

The Chernoff bound states that for independent random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$, letting $\mu = \mathbb{E}\left[\sum_{i=1}^{n} X_i\right]$, for any $\delta > 0$,

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

2. Assume there are 1000 registered users on your site $u_1, \ldots, u_{1000}$, and in a given day, each user visits the site with some probability $p_i$. The event that any user visits the site is independent of what the other users do. Assume that $\sum_{i=1}^{1000} p_i = 500$.

   (a) Let **X** be the number of users that visit the site on the given day. What is $\mathbb{E}[\mathbf{X}]$.
   (b) Apply a Chernoff bound to show that $\Pr[\mathbf{X} \geq 600] \leq .01$. $1/5$
   (c) Apply Markov's inequality and Chebyshev's inequality to bound the same probability. How do they compare?

a) $X = \sum X_i \rightarrow \mathbb{E}X = \sum_{i=1}^{1000} \mathbb{E}X_i = \sum_{i=1}^{1000} p_i = 500$

b) $\Pr\left(\sum X_i - 500 \geq \delta 500\right) \leq 2\exp\left(\frac{\delta^2 \mu}{2\delta}\right)$    $\boxed{\delta = 1/5}$

$\Pr\left(\sum X_i - 500 \geq 100\right) \leq 2\exp\left(\frac{1/5^2 \cdot 500}{2 + 1/5}\right)$

$2\exp\left(\frac{20}{4}\right) \leq 2\exp(-5) \leq .01$

The Chernoff bound states that for independent random variables $X_1, \ldots, X_n$ taking values in $\{0, 1\}$, letting $\mu = \mathbb{E}\left[\sum_{i=1}^{n} X_i\right]$, for any $\delta > 0$,
$$\Pr\left(\left|\sum_{i=1}^{n} X_i - \mu\right| > \delta\mu\right) \leq 2\exp\left(-\frac{\delta^2\mu}{2+\delta}\right).$$

11

ALWAYS, SOMETIMES, or NEVER:

$$\left(1 - \underline{\delta}\right)^n = 1 - \delta n + \delta^2 \ldots + \delta^?$$

2. $\Pr[\max(X_1, \ldots X_n) \underset{=}{\geq} t] \leq \sum_{i=1}^{n} \Pr[X_i \geq t]$ for any random variables $X_1, \ldots, X_n$.

$$\Pr\left(X_1 \geq t \text{ OR } X_2 \geq t \text{ OR } \ldots X_n \geq t\right) \leq \sum \Pr\left(X_i \geq t\right)$$

union bound

HW which answers middle is
use this to answer $n$ questions    with    prob $1-\delta$

(c) $\Pr[\mathbf{X} = s \cap \mathbf{Y} = t] = \Pr[\mathbf{X} = s] \cdot \Pr[\mathbf{Y} = t].$  ✓ independent

$$\Pr\left(\text{get all questions correct}\right) \geq 1 - n\delta$$

$$1 - \Pr\left(\text{fails at least one}\right)$$
union bound    $\leq n \cdot \delta = \sum_{i=1}^{n} \Pr\left(\text{fails question } i\right)$