# CS-GY 9223I (CS-UY 3943)

# Algorithmic Machine Learning + Data Science

**Instructor:** Prof. Christopher Musco (cmusco@nyu.edu)

**Office Hours:** 3.5pm Thursdays (except tomorrow)

**Reading Group:** TBD. Poll to gauge interest and time.

**Webpage:** chrismusco.com/9223_2019

NYU Classes for forum

**Exams:** Midterm - 10/23, in class → check for conflicts now!

Final - 12/18, class time

# Course Topic: Algorithmic methods for machine learning + data analysis at scale.

- High throughput / realtime data applications

(think Shazam, Google maps/waze, Amazon product recs, industrial robotics, FinTech, scientific applications)

- More complex models → more training data

(deep neural networks, reinforcement learning, machine translation)

- Data analysis on low compute devices

(smart phones/watches, robots/drones/etc., sensor networks)

Some numbers:

- Google receives $\approx$ 10,000 Maps queries <u>every second</u>

- NASA collects 6.4 Tb of satellite images <u>every day</u>

  - new "ImageNet" dataset every 3 days

- Large Synoptic Survey Telescope : 15 Tb of images
  <u>per night</u>

- Broad Institute sequences 24 Tb
  of genetic data per day

"needle in a haystack problems"

Ushering in new golden age for research in computational methods, using entirely new tools!

# Course Objectives:

**1** New algorithmic toolkit (randomization, sketching, optimization, spectral methods, etc.)

  - lectures    - reading

  See   cs.princeton.edu/courses/archive/fall18/cos521

**2.** Learn to apply tools in the world (industry, academia)

  - 4 Problem Sets    - in class work

$\longrightarrow$ midway into class

$\longrightarrow$ break into groups (auditors included)

$\longrightarrow$ pset like problem to solve

**3.** Theory as an approach to algorithm design.

# What we won't cover:

- Software tools or frameworks
  (MapReduce, Tensorflow, Amazon AWS)

  (CS-GY 6513: Big Data)

- Machine learning models
  (neural nets, RL, Bayesian methods, unsupervised
  learning, function fitting, etc.)

# Unit 1: The Power of Randomness

## Hashing (+ Load Balancing)

- Work horse of the modern web    - good probability review!

## Probability review:

X takes values in set $S \subseteq \mathbb{R}$. Eg. $S = \{1, 2, 3, 4, 5, 6\}$ for dice roll.

Expectation: $\mathbb{E}[X] = \sum_{\omega \in S} Pr[X = \omega] \cdot \omega$

Continuous r.v. s : $\mathbb{E}[X] = \int_{y \in \mathbb{R}} p(y) \, dy$

## Independence

Two random events $A, B$ are independent if $Pr(A|B) = Pr(A)$

$\underset{\nwarrow}{}$

"A given B"

$$Pr(A|B) \overset{def}{=} \frac{Pr(A \cap B)}{P(B)}$$

So equivalently, when $A$ and $B$ are independent,

$$\frac{Pr(A \cap B)}{Pr(B)} = Pr(A) \longrightarrow Pr(A \cap B) = Pr(A) \cdot Pr(B)$$

Roll 2 dice. What's the probability the first is odd and the second is < 3?

Independence   Given random variables X and Y taking
values in $S_x$ and $S_y$. We say X and Y are
independent if for all $w \in S_x$, $z \in S_y$
   $[X = w]$ and $[Y = z]$ are independent random
                                                events.

# Expectation Identities

$$\mathbb{E}[aX] = a\mathbb{E}[X]$$

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y] \; ?$$

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \; ?$$

$\mathbb{E}[X+Y]$ is true for <u>any</u> $X, Y$.

$$= \sum_{w \in S_x} \sum_{z \in S_y} \Pr(X=w, Y=z) \cdot (w+z)$$

$$= \sum_w \sum_z \Pr(X=w, Y=z) \cdot w + \sum_w \sum_z \Pr(X=w, Y=z) \cdot z$$

$$= \sum_w w \cdot \underbrace{\sum_z \Pr(X=w, Y=z)}_{= \Pr(X=w)} + \sum_z z \cdot \underbrace{\sum_w \Pr(X=w, Y=z)}_{= \Pr(Y=z)}$$

$$= \mathbb{E}[X] + \mathbb{E}[Y]$$

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \quad \text{true for } \underline{\text{independent}} \text{ r.v.'s.}$$

$$\mathbb{E}[XY] = \sum_{w \in S_X} \sum_{z \in S_Y} \Pr(X=w, Y=z) \, w \, z$$

$$= \sum_w \sum_z \Pr(X=w) \Pr(Y=z) \, w \cdot z$$

$$= \sum_w \left[ w \cdot \Pr(X=w) \cdot \sum_z \Pr(Y=z) \cdot z \right]$$

$$= \left[ \sum_w w \cdot \Pr(X=w) \right] \cdot \left[ \sum_z \Pr(Y=z) \cdot z \right]$$

$$= \mathbb{E}[X] \cdot \mathbb{E}[Y]$$

$$\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y] \implies X, Y \text{ are "uncorrelated"}.$$

$$\text{Independence} \implies \text{Uncorrelated}. \quad \text{Uncorrelated} \nRightarrow \text{Independence}.$$

# Markov's [least exciting] Inequality

For a non-negative random variable $X$,

$$\Pr[X \geq a] \leq \mathbb{E}[X]/a$$

Equivalent: $\Pr[X \geq c \cdot \mathbb{E}[X]] \leq 1/c$.   Think $c = 2, 10, \ldots$

"concentration inequality"

Proof:

$$\mathbb{E}[X] = \sum_{\omega} \Pr[X = \omega] \cdot \omega$$

$$= \underbrace{\sum_{\omega < a} \Pr[X = \omega] \cdot \omega}_{\geq 0} + \underbrace{\sum_{\omega \geq a} \Pr[X = \omega] \cdot \omega}$$

$$\geq \sum_{\omega \geq a} \Pr[X = \omega] \cdot a$$

$$= a \cdot \sum_{\omega \geq a} \Pr[X = \omega]$$

$$= a \cdot \Pr[X \geq a]$$

$$\mathbb{E}[X] \geq a \cdot \Pr[X \geq a]$$

# Is Markov's Inequality Tight? Can you prove something better?

In general, Markov's is tight.

$X = 0$ with probability $1 - \frac{t}{a}$

$\phantom{X} = a$ with probability $t/a$.

$$\mathbb{E}[x] = 0 \cdot \left(1 - \frac{t}{a}\right) + a \cdot (t/a) = t.$$

$$\mathbb{P}[x \geq a] = t/a = \mathbb{E}[x]/a.$$

# Hashing (+ Load Balancing)

- Work horse of the modern web
- good probability review!

(key, value) store

3 operations:

"nyu.edu"
insert $(K_1, v_1)$ → 216.165.47.10
insert $(K_2, v_2)$
$\vdots$
insert $(K_m, v_m)$

$\vdots$

delete $(K_{10}, v_{10})$
$\vdots$

query $(k_j)$ → $v_j$
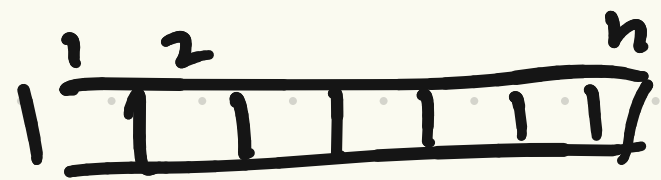query $(k_{10})$ → empty

↓

Want

1) Fast queries:
   $O(1)$ time

2) Small space:
   $O(m)$ space

# Hashing

- Build table $T$ of size $n$



- Choose __random__ function $\quad h: U \to \{1, \ldots, n\}$

$\downarrow$

universe of possible keys

$k_1, \ldots, k_m \in U$

| $U$ | $h$ |
|---|---|
| $u_1$ | $h(u_1) = 10$ |
| $u_2$ | $h(u_2) = 4$ |
| $u_3$ | $h(u_3) = m$ |
| $\vdots$ | $\vdots$ |
| | $h(u_{100}) = 4$ |
| $u_{1000000}$ | |

$h$ drawn uniform from $\mathcal{H}$

all possible mappings from

$U \to \{1, \ldots, n\}$

"hash family"

$(k \cdot r_1 + r_2)(\text{mod } m)$

Is this possible in practice?

3 Issues!

# Hashing

- for insert$(k,v)$, store $v$ at $T_{h(k)}$
- for delete$(k,v)$, remove $v$ from $T_{h(k)}$
- for query$(k)$, look at $T_{h(k)}$

$m \ll |U|$ so could have $h(k) = h(j)$ for $j \neq k$.

"hash collision"
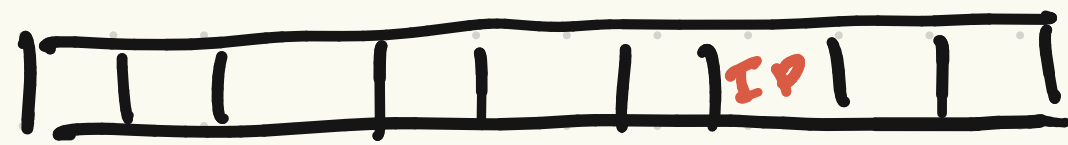
Store $v(k), v(j)$ in linked list at $T_{h(k)} = T_{h(j)}$

Goal: Hash collisions should be rare!

Lookup time $= O(1)$ if no collision,

otherwise $O(\text{length of linked list})$

# Hashing Applications

- URL / DNS resolution.   What IP address to visit
  to see www.nyu.edu?



$$T_h("www.nyu.edu")$$

- Web content delivery (distributed hash table)

- Amazon Dynamo DB, Mongo DB, Cassandra, Google
  Data store
  "no-sql data bases

- Google directions: send "boston to NYC" to   $A_{h("boston to NYC")}$
  servers $A_1, \ldots, A_8$

# Goal: Hash collisions should be rare!

How many collisions in _expectation_ when inserting $\underline{m}$ items into table of size $\underline{n}$?

#### of collisions = random variable, with randomness coming from choice of $h \in \mathcal{H}$.

take $K_1, \ldots, K_m$ as fixed.

$$C := \frac{1}{2} \sum_{i=1,\ldots,m} \sum_{j \neq i} \mathbb{1}[h(K_i) = h(K_j)]$$

$\mathbb{1}[\text{true}] = 1$
$\mathbb{1}[\text{false}] = 0$

"indicator function"

How many collisions in _expectation_ when
inserting $\underline{m}$ items into table of size $\underline{n}$ ?

$$\mathbb{E}[c] = \mathbb{E}\left[\frac{1}{2} \sum_{i} \sum_{j \neq i} \mathbb{1}[h(k_i) = h(k_j)]\right]$$

$$= \frac{1}{2} \sum_{i} \sum_{j \neq i} \underbrace{\mathbb{E}[\mathbb{1}[h(k_i) = h(k_j)]]}_{= 1/n}$$

$$= \frac{1}{2} \sum_{i} \sum_{j \neq i} 1/n \quad = \boxed{\frac{m \cdot m - 1}{2n}}$$

$\longrightarrow$ O(1) time lookups.

**Result 1**  **Collision free** hash table with $O(m^2)$ space.

Set $n = 5 m^2$.  $E[C] = \frac{m(m-1)}{2n} \leq \frac{1}{10}$.  $\underline{Pr[C \geq 1] \leq \frac{1}{10}}.$

Markov's inequality

Could keep retrying until achieve collision free hash.

Trials:          1          2          3          4    . . . .        6

Failure Probability:  $\frac{1}{10}$    $\frac{1}{100}$    $\frac{1}{1000}$    $\frac{1}{10000}$         $\leq$ chance getting struck by lightening.
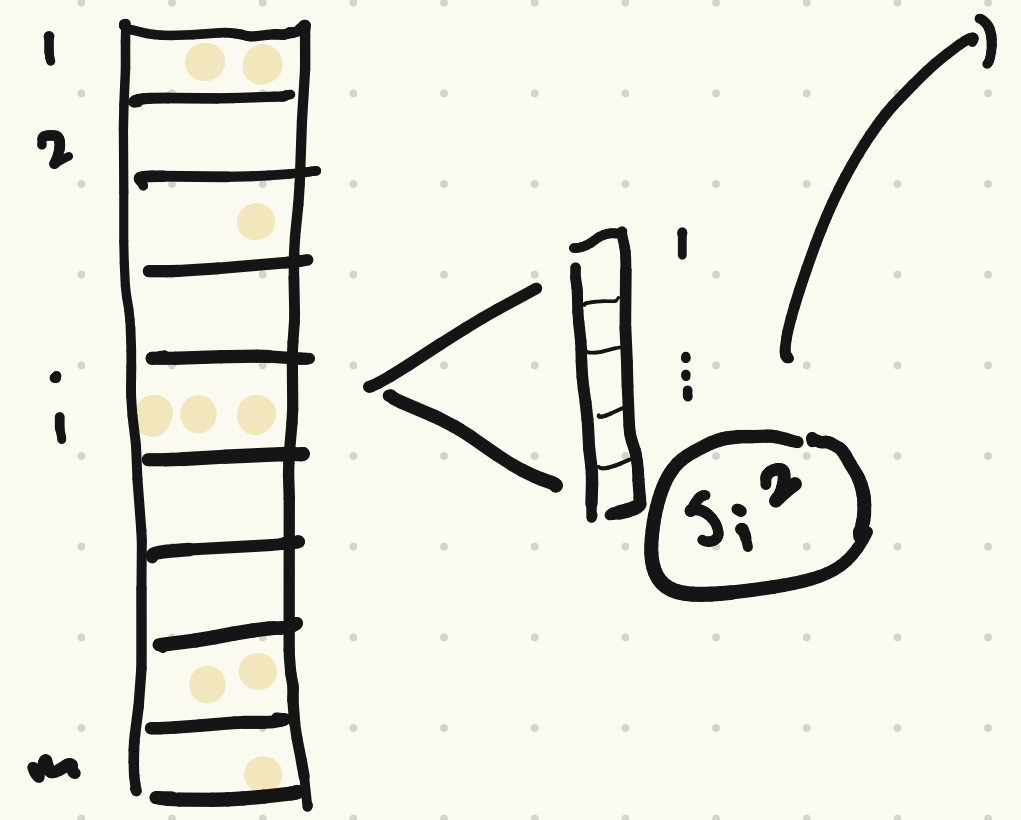
$O(m^2)$ is a lot of space overhead for $m$ items.

"Birthday Paradox"

# Result 2   Collision free hash table with $O(m)$ space.

## Key Idea :   2 Layer Hashing

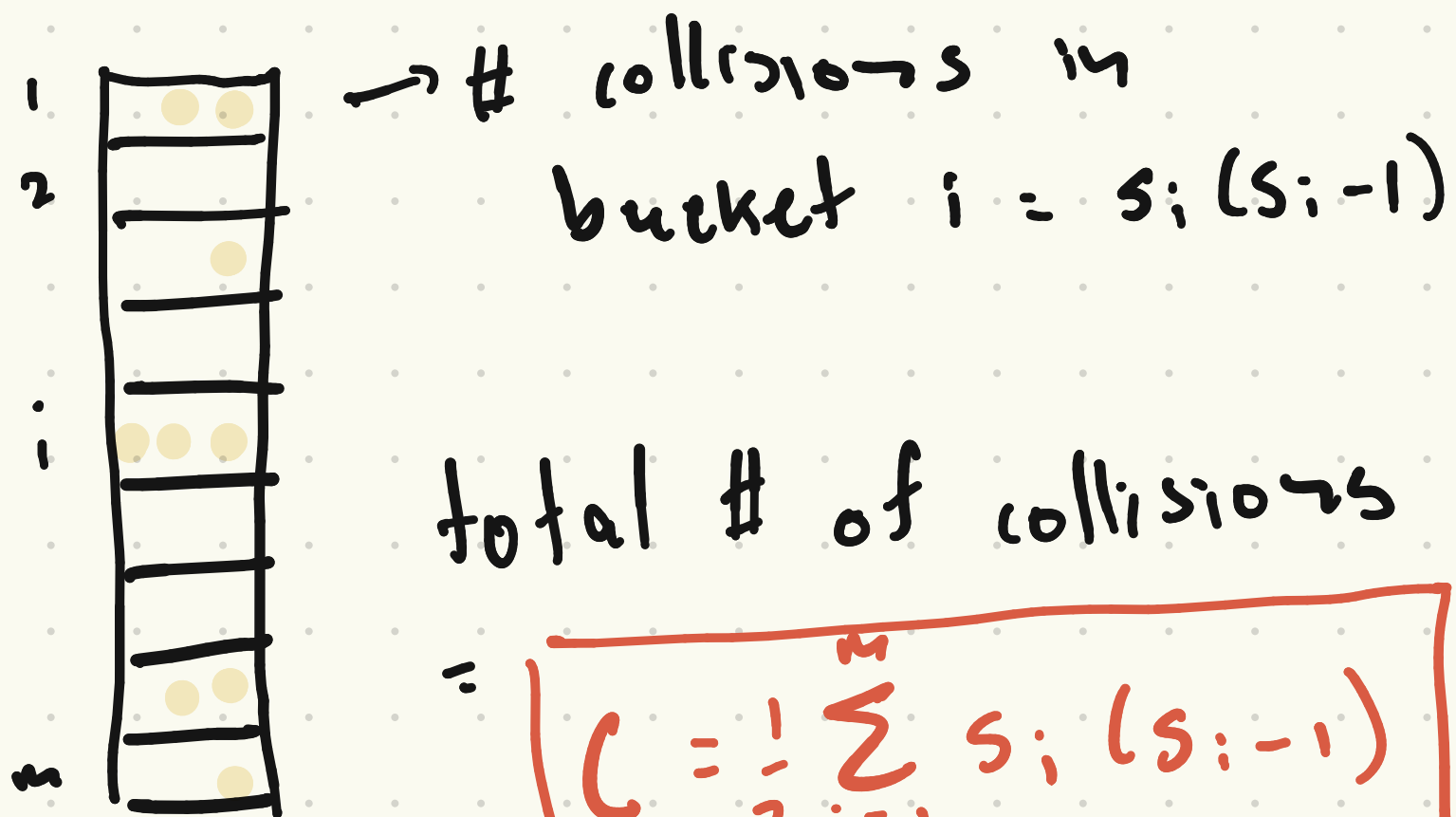2nd level is collision free $\to$ each lookup takes just 2 hash function evaluations.

$$\text{Total space} = m + \sum_{i=1}^{m} s_i^2$$

$$\mathbb{E}\left[ m + \sum_{i=1}^{m} s_i^2 \right] = m + \mathbb{E} \sum_{i=1}^{m} s_i^2$$

$\underbrace{\phantom{\mathbb{E} \sum_{i=1}^{m} s_i^2}}_{?}$

# items in bin $i = s_i$

$$\mathbb{E}[\text{total space}] = m + \mathbb{E}\left[\sum_{i=1}^{m} s_i^2\right]$$



$\rightarrow$ # collisions in bucket $i = s_i(s_i - 1)$

total # of collisions

$= \boxed{C = \frac{1}{2}\sum_{i=1}^{m} s_i(s_i - 1)}$ $\longrightarrow$ almost $s_i^2$!

Almost instant lookups with only $3\times$ space overhead!

$\rightarrow$ Markov's inequality is actually quite powerful!

$$\mathbb{E}\sum_{i=1}^{m} s_i^2 = \mathbb{E}\left[\sum_{i=1}^{m} s_i(s_i-1) + \sum_{i=1}^{m} s_i\right]$$

$$= 2\mathbb{E}[C] + m \quad < 2m$$

with $m$ buckets, $m$ items
$= \frac{m(m-1)}{2m} \leq \frac{1}{2}$

$$\boxed{\mathbb{E}[\text{total space}] = 3m}$$

- Google directions: send "boston to NYC" to $A_{h("boston\ to\ NYC")}$

$A_1, A_2, \ldots, A_n$

why ?

– Google directions: send "boston to NYC" to $A_{h("boston \; to \; NYC")}$

$A_1, A_2, \ldots, A_n$

Now we care about the maximum # of elements per bin.

"Load Balancing" $\max[S_1, \ldots, S_n]$

Suppose we hash $\underline{n \; items}$ to $\underline{n \; servers}$. "Balls into Bins"

$\rightarrow b_1, \ldots, b_n$

What does Markov give?

$$\mathbb{E}[S_i] = \mathbb{E}\left[\sum_{j=1}^{n} \mathbb{1}[h(b_j) = i]\right] = \sum_{j=1}^{n} \frac{1}{n} = 1.$$

$$\boxed{\Pr[S_i \geq 10] \leq \frac{1}{10}}$$

## 1 in 10 servers could be overloaded! No bound on max.

**Goal**: For <u>any</u> $i$, $\Pr[s_i \geqslant B] \leq \frac{1}{10n}$.

**Corollary**: With probability $9/10$ <u>all</u> $s_i \leq B$.

**Proof**: <span style="color:red">Union bound: Given random events $A$, $B$

$\Pr[A \text{ or } B \text{ occur}] \leq \Pr[A] + \Pr[B]$

$= \Pr[A + \bar{B}] + \Pr[\bar{A} + B] + \Pr[A + B]$

$= \Pr[A + \bar{B}] + \Pr[B] \leq \Pr[A] + \Pr[B]$</span>

$\Pr[s_1 \geqslant B \text{ or } s_2 \geqslant B \text{ or} \ldots s_n \geqslant B] \leq \frac{1}{10n} + \frac{1}{10n} + \ldots + \frac{1}{10n} = \frac{1}{10}$.

**Goal**: For <u>any</u> $i$, $\Pr[s_i \geq B] \leq \frac{1}{10n}$.

Markovs: $B = 10n$. $\longrightarrow$ Vacuous bound! Only $n$ items...

Tighter bounds on <u>deviation</u> by considering more information about a random variable.

<u>Variance</u>: For r.v. $X$, $\mathrm{Var}(X) = \mathbb{E}[(X - \mathbb{E}X)^2]$

Equivalent: $\mathrm{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$

<u>Proof</u>: $\mathbb{E}[(X - \mathbb{E}X)^2] = \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2]$

$$Var(\alpha X) = \alpha^2 Var(x)$$

Linearity of
Variance: $Var[X+Y] = Var[X] + Var[Y]$ if $X, Y$ independent

$Var[X_1 + X_2 + \ldots + X_m] = Var[X_1] + Var[X_2] + \ldots + Var[X_m]$

if $X_i, X_j$ independent for all $i, j \in 1, \ldots, m$.

"pairwise independence"

"mutual independence" is stronger. why?

# Chebyshev Inequality: Let $X$ be a r.v. with $\text{var}(X) = \sigma^2$.

$$\Pr[|X - \mathbb{E}X| \geqslant k\sigma] \leqslant 1/k^2$$

upper and lower bound          — no assumption that $X > 0$.

Proof. Consider the random variable $Y = (X - \mathbb{E}X)^2$.

$Y$ is non negative. By Markov inequality,

$$\Pr[Y \geqslant k^2 \mathbb{E}[Y]] \leqslant 1/k^2$$

$$\mathbb{E}[Y] = \mathbb{E}[X - \mathbb{E}X] = \text{var}(X) = \sigma^2.$$

$$\Pr[(X - \mathbb{E}X)^2 \geqslant k^2 \sigma^2] \leqslant 1/k^2$$

$$\boxed{\text{Goal: For any } i, \quad \Pr[S_i \geq B] \leq \frac{1}{10n}}$$

$B$ = # of balls to bin $i$

$$S_i = \sum_{j=1}^{n} \mathbb{1}[h(b_j) = i] \qquad \text{Var}[S_i] = \sum_{j=1}^{n} \text{Var} \underbrace{\mathbb{1}[h(b_j) = i]}_{X_j}$$

$$\text{Var}(X_j) = \mathbb{E}[X_j^2] - \mathbb{E}[X_j]^2 = \boxed{\frac{1}{n} - \frac{1}{n^2}}$$

$$X_j = \begin{cases} 1 & \text{w/ prob. } 1/n \\ 0 & \text{otherwise} \end{cases} \qquad X_j^2 = \begin{cases} 1 & \text{w/ prob } 1/n \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{E}[X_j] = 1/n \qquad\qquad \mathbb{E}[X_j^2] = 1/n$$

$$\text{Var}[S_i] = n\left(\frac{1}{n} - \frac{1}{n^2}\right) = 1 - \frac{1}{n} \leq 1.$$

**Goal**: For any $i$, $\Pr[s_i \geq B] \leq \frac{1}{10n}$.

Chebyshev's: $\Pr[|x - \mathbb{E}x| \geq k\sigma] \leq 1/k^2$

$$\text{Var}[s_i] = 1 - 1/n.$$

$$\mathbb{E}[s_i] = 1$$

Improvable to $O(\log n)$!

Set $k = \sqrt{10n}$.

$$\Pr\left[|s_i - 1| \geq \sqrt{10n} \cdot \sqrt{1 - 1/n}\right] \leq \frac{1}{10n}$$

$$\Pr\left[s_i - 1 \geq \sqrt{10n} \cdot 1\right] \leq 1/10\,]$$

With probability $9/10$ <u>all</u> bins have $\leq O(\sqrt{n})$ balls.

$$\boxed{\Pr[s_i \geq O(\sqrt{n})] \leq 1/10n}$$

$n = 1,000,000$

max load $\approx 1000$.

# In class exercise

New York University Tandon School of Engineering
Computer Science and Engineering

## CS-GY 9223I: Lecture 1 Coursework

### Problem 1: Hash collisions are useful?

Your company is considering paying for a cloud service that provides CAPTCHA-like visual puzzles for verifying that users are human. The company providing the service claims to have a larger database of unique puzzles than any competitors, but you don't trust the salesperson.

(a) The company has provided you with an API end-point which returns puzzles uniformly and independently at random from their database. Using this endpoint, describe a simple randomized estimator for the number of puzzles in the database, $n$.

(b) The company claims their database has $1,000,000$ unique CAPTCHAs in it. Using your estimator, roughly how many queries do you need to verify their claim with good probability (e.g. 9/10)? You should need far less than $1,000,000$ queries!

(c) More generally, how many samples are required to estimate the true number of CAPTCHAs, $n$, in the database up to additive error $\pm \epsilon n$, with good probability?