

## CS-GY 6763: Lecture 3

# Finish Chebyshev's, Exponential Concentration Inequalities

---

NYU Tandon School of Engineering, Prof. Christopher Musco

## DISTINCT ELEMENTS PROBLEM

**Input:**  $x_1, \dots, x_n \in \mathcal{U}$  where  $\mathcal{U}$  is a huge universe of items.

**Output:** Number of distinct inputs,  $D$ .

**Example:**  $f(\underline{1}, \underline{10}, \underline{2}, \underline{4}, \underline{9}, \underline{2}, \underline{10}, \underline{4}) \rightarrow D = 5$

---

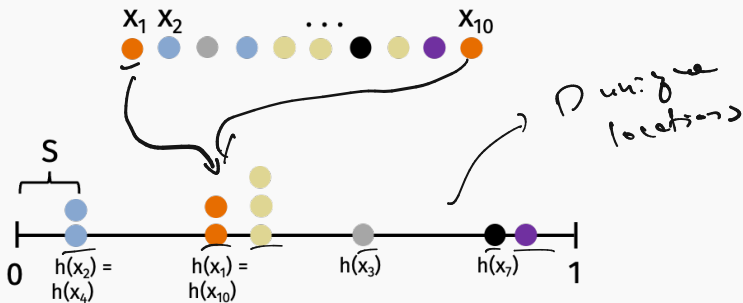
Flajolet-Martin (simplified):

- Choose random hash function  $h : \mathcal{U} \rightarrow [0, 1]$ .
- $S = 1$
- For  $i = 1, \dots, n$ 
  - $S \leftarrow \min(S, h(x_i))$
- Return  $\frac{1}{S} - 1$

[ 0 1 ]

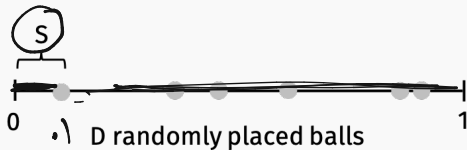
## Flajolet–Martin (simplified):

- Choose random hash function  $h : \mathcal{U} \rightarrow [0, 1]$ .
- $S = 1$
- For  $i = 1, \dots, n$ 
  - $S \leftarrow \min(S, h(x_i))$
- Return:  $\tilde{D} = \frac{1}{S} - 1$



## PROOF "FROM THE BOOK"

$$\boxed{\mathbb{E}[S]} = \Pr[(D+1)^{\text{st}} \text{ item has the smallest hash value}] = \frac{1}{D+1}$$



By symmetry, this equals  $\frac{1}{D+1}$  (since every ball is equally likely to be first).

Final Estimate:  $\tilde{D} = \frac{1}{S} - 1$

$$D+1 = \frac{1}{\mathbb{E}[S]}$$

$$D = \frac{1}{\mathbb{E}[S]} - 1$$

## PROVING CONCENTRATION

$\mathbb{E}S = \frac{1}{D+1}$ . Estimate:  $\tilde{D} = \frac{1}{S} - 1$ . Claim: We have for  $\epsilon < \frac{1}{2}$ :

If  $(1 - \epsilon)\mathbb{E}S \leq S \leq (1 + \epsilon)\mathbb{E}S$ , then:

$$(1 - 4\epsilon)\tilde{D} \leq \underline{\tilde{D}} \leq (1 + 4\epsilon)\tilde{D}.$$

$$\frac{1}{(1-\epsilon)} \frac{1}{\mathbb{E}[S]} \geq \frac{1}{S} \geq \frac{1}{1+\epsilon} \frac{1}{\mathbb{E}[S]} \Rightarrow \underline{(D+1) \frac{1}{1+\epsilon}} \leq \frac{1}{S} \leq (D+1) \frac{1}{1-\epsilon}$$

$$(1-\epsilon)(D+1) \leq \frac{1}{S} \leq (1+2\epsilon)(D+1) \Rightarrow (1-\epsilon)D + 1 - \epsilon \leq \frac{1}{S} \leq (1+2\epsilon)D + 1 + 2\epsilon$$

$$(1-\epsilon)D - \epsilon \leq \frac{1}{S} - 1 \leq (1+2\epsilon)D + 2\epsilon$$

$$\boxed{(1-2\epsilon)D \leq \frac{1}{S} - 1 \leq (1+4\epsilon)D}$$

So, it suffices to show that  $S$  concentrates around its mean. I.e. that  $|S - \mathbb{E}S| \leq \epsilon \cdot \mathbb{E}S$ . We will use Chebyshev's inequality as our concentration bound.

Recall:

$$1 + \epsilon \leq \frac{1}{1 - \epsilon} \leq 1 + 2\epsilon \text{ for } \epsilon \in [0, .5].$$

$$\underline{1 - \epsilon} \leq \underline{\underline{\frac{1}{1 + \epsilon}}} \leq 1 - .5\epsilon \text{ for } \epsilon \in [0, 1].$$

## CALCULUS PROOF

Lemma

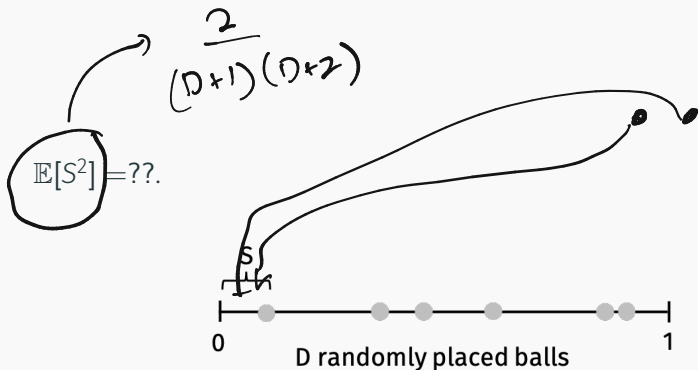
$$\underline{\text{Var}[S]} = \underline{\mathbb{E}[S^2]} - \underline{\mathbb{E}[S]^2} = \frac{2}{(D+1)(D+2)} - \frac{1}{(D+1)^2} \leq \frac{1}{(D+1)^2}.$$

Proof:

$$\begin{aligned}\underline{\mathbb{E}[S^2]} &= \int_0^1 \underbrace{\text{Pr}[S^2 \geq \lambda]}_{\substack{S \\ \geq \sqrt{\lambda}}} d\lambda \\ &= \int_0^1 \underbrace{\text{Pr}[S \geq \sqrt{\lambda}]}_{\substack{S \\ \geq \sqrt{\lambda}}} d\lambda \\ &= \int_0^1 \underbrace{(1 - \sqrt{\lambda})^D}_{\substack{S \\ \geq \sqrt{\lambda}}} d\lambda \\ &= \frac{2}{(D+1)(D+2)}\end{aligned}$$

[www.wolframalpha.com/input?i=antiderivative+of+%281-sqrt%28x%29%29%5ED](http://www.wolframalpha.com/input?i=antiderivative+of+%281-sqrt%28x%29%29%5ED)

# PROOF "FROM THE BOOK"



$E[s^2] = \Pr[(D+1)\text{st and } (D+2)\text{nd balls are two squares apart}].$

$$\Downarrow \frac{2}{(D+2)(D+1)}$$



Recall we want to show that, with high probability,

$$(1 - \epsilon)\mathbb{E}[S] \leq \underline{S} \leq (1 + \epsilon)\mathbb{E}[S].$$

$$\text{Var}[S] = \frac{2}{(D+2)(D+1)} - \frac{1}{(D+1)^2}$$

$$\leq \frac{2}{(D+1)^2} - \frac{1}{(D+1)^2}$$

$$\cdot \mathbb{E}[S] = \frac{1}{D+1} = \mu.$$

$$\cdot \text{Var}[S] \leq \frac{1}{(D+1)^2} = \mu^2 \quad \text{Standard deviation: } \underline{\sigma \leq \mu}.$$

$$\cdot \text{Want to bound } \Pr[|S - \mu| \geq \epsilon\mu] \leq \delta.$$

$$\text{Chebyshev's: } \Pr[|S - \mu| \geq \epsilon\mu] = \Pr[|S - \mu| \geq \epsilon\sigma] \leq \frac{1}{\epsilon^2}.$$

**Vacuous bound. Our variance is way too high!**

$$\Pr[|S - \mu| \geq \underbrace{k\mu}_{\epsilon\mu}] \leq \frac{1}{k^2} \quad k = \epsilon$$

## VARIANCE REDUCTION

$$X_i \sim \mu$$

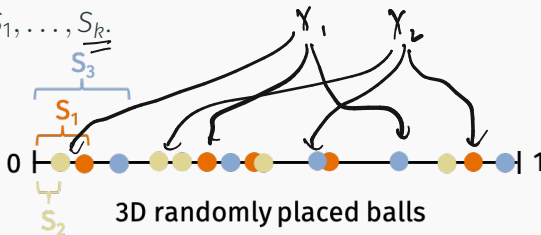
**Trick of the trade:** Repeat many independent trials and take the mean to get a better estimator.

Given i.i.d. (independent, identically distributed) random variables  $X_1, \dots, X_k$  with mean  $\mu$  and variance  $\sigma^2$  what is:

$$\bullet \mathbb{E} \left[ \frac{1}{k} \sum_{i=1}^k X_i \right] = \frac{1}{k} \sum_{i=1}^k \mathbb{E}[X_i] = \frac{1}{k} \sum_{i=1}^k \mu = \frac{1}{k} \cdot k \mu = \underline{\underline{\mu}}$$

$$\begin{aligned} \bullet \text{Var} \left[ \frac{1}{k} \sum_{i=1}^k X_i \right] &= \frac{1}{k^2} \text{Var} \left[ \sum_{i=1}^k X_i \right] = \frac{1}{k^2} \sum_{i=1}^k \underbrace{\text{Var}[X_i]}_{\sigma^2} \\ &= \frac{1}{k^2} \cdot k \cdot \sigma^2 = \left( \frac{1}{k} \cdot \sigma^2 \right) \end{aligned}$$

Using independent hash functions, maintain  $k$  independent sketches  $S_1, \dots, S_k$ .



Flajolet-Martin:

- Choose  $k$  random hash functions  $h_1, \dots, h_k : \mathcal{U} \rightarrow [0, 1]$ .
- $S_1 = 1, \dots, S_k = 1$
- For  $i = 1, \dots, n$ 
  - $S_j \leftarrow \min(S_j, h_j(x_i))$  for all  $j \in 1, \dots, k$ .
- $S = (S_1 + \dots + S_k) / k$
- Return:  $\frac{1}{S} - 1$

1 estimator:

$$\mathbb{E}[S] = \frac{1}{D+1} = \mu.$$

$$\text{Var}[S] \leq \mu^2$$

(Chebyshev's:

$$\Pr[|S - \mathbb{E}[S]| \geq c \frac{\mu}{\sqrt{k}}] \leq \frac{1}{c^2}$$

$$\underbrace{c}_{\epsilon \mu}$$

$$\frac{c}{\sqrt{k}} = \epsilon$$

$$k = \frac{c^2}{\epsilon^2}$$

 $k$  estimators:

$$\mathbb{E}[S] = \frac{1}{D+1} = \mu.$$

$$\text{Var}[S] \leq \mu^2/k$$

$$\text{By Chebyshev, } \Pr[|S - \mathbb{E}[S]| \geq c\mu/\sqrt{k}] \leq \frac{1}{c^2}.$$

$$\rightarrow = \delta$$

Setting  $c = 1/\sqrt{\delta}$  and  $k = \frac{1}{\epsilon^2 \delta}$  gives:

$$k = \frac{c^2}{\epsilon^2} \rightarrow 1/\delta$$

$$k = \frac{1}{\delta \epsilon^2}$$

$$\Pr[|S - \mu| \geq \epsilon \mu] \leq \delta.$$

- Recall that to ensure  $(1 - \bar{\epsilon})D \leq \frac{1}{S} - 1 \leq (1 + \bar{\epsilon})D$ , we needed  $|\underline{S} - \mu| \leq \frac{\bar{\epsilon}}{4}\mu$ .
- So apply the result from the previous slide with  $\epsilon = \bar{\epsilon}/4$ .
- Need to store  $k = \frac{1}{\epsilon^2 \delta} = \frac{1}{(\bar{\epsilon}/4)^2 \delta} = \frac{16}{\bar{\epsilon}^2 \delta}$  counters.

**Total space complexity:**  $O\left(\frac{1}{\epsilon^2 \delta}\right)$  to estimate distinct elements up to error  $\epsilon$  with success probability  $1 - \delta$ .

## NOTE ON FAILURE PROBABILITY

$$\log(1/\delta)$$

$O\left(\frac{1}{\epsilon^2 \delta}\right)$  space is an impressive bound:

- $1/\epsilon^2$  dependence cannot be improved.
- No linear dependence on number of distinct elements  $D$ .<sup>1</sup>
- But...  $1/\delta$  dependence is not ideal. For 95% success rate, pay a  $\frac{1}{5\%} = \underline{20}$  factor overhead in space.

We can get a better bound depending on  $O(\log(1/\delta))$  using exponential tail bounds. We will see next.

---

<sup>1</sup>Technically, if we account for the bit complexity of storing  $S_1, \dots, S_k$  and the hash functions  $h_1, \dots, h_k$ , the space complexity is  $O\left(\frac{\log D}{\epsilon^2 \delta}\right)$ .

## DISTINCT ELEMENTS IN PRACTICE

In practice, we cannot hash to real numbers on  $[0, 1]$ . Could use a finite grid, but more popular choice is to hash to integers (bit vectors).

### Real Flajolet-Martin / HyperLogLog:

$h(x_1)$	101001 <u>0</u>
$h(x_2)$	10011 <u>00</u>
$h(x_3)$	10011 <u>0</u>
⋮	
$h(x_n)$	1011 <u>000</u>

- Estimate # distinct elements based on maximum number of trailing zeros  $m$ .
- The more distinct hashes we see, the higher we expect this maximum to be.

## LOGLOG SPACE

Total Space:  $O\left(\frac{\log \log D}{\epsilon^2} + \log D\right)$  for an  $\epsilon$  approximate count.

“Using an auxiliary memory smaller than the size of this abstract, the LogLog algorithm makes it possible to estimate in a single pass and within a few percents the number of different words in the whole of Shakespeare’s works.” – Flajolet, Durand.

$$\frac{\log D}{\epsilon^2}$$



**Total Space:**  $O\left(\frac{\log \log D}{\epsilon^2} + \log D\right)$  for an  $\epsilon$  approximate count.

“Using an auxiliary memory smaller than the size of this abstract, the LogLog algorithm makes it possible to estimate in a single pass and within a few percents the number of different words in the whole of Shakespeare’s works.” – Flajolet, Durand.

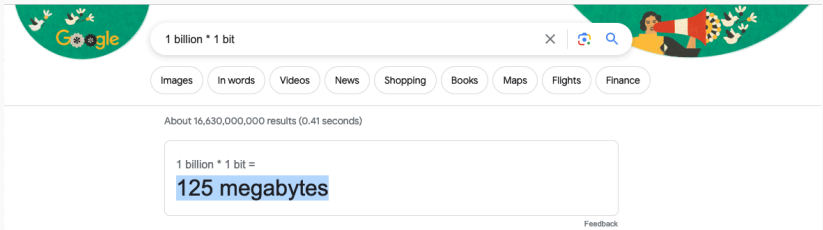
**Total Space:**  $O\left(\frac{\log \log D}{\epsilon^2} + \log D\right)$  for an  $\epsilon$  approximate count.

“Using an auxiliary memory smaller than the size of this abstract, the LogLog algorithm makes it possible to estimate in a single pass and within a few percents the number of different words in the whole of Shakespeare’s works.” – Flajolet, Durand.

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\begin{aligned}
 \text{space used} &= O\left(\frac{\log \log D}{\epsilon^2} + \log D\right) \\
 &= \frac{1.04 \cdot \lceil \log_2 \log_2 D \rceil}{\epsilon^2} + \lceil \log_2 D \rceil \text{ bits} \\
 &= \frac{1.04 \cdot 5}{.02^2} + 30 = \underline{13030} \text{ bits} \approx \underline{1.6 \text{ kB!}}
 \end{aligned}$$

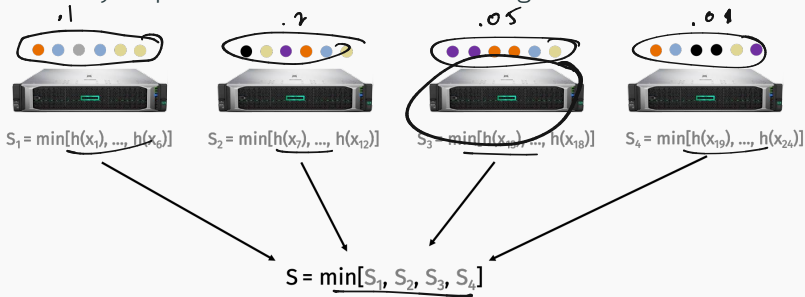
Although, to be fair, storing a dictionary with (1 billion) bits only takes 125 megabytes. Not tiny, but not unreasonable.



These estimators become more important when you want to count many different things (e.g., a software company tracking clicks on 100s of UI elements).

## DISTRIBUTED DISTINCT ELEMENTS

Also very important in distributed settings.



Distinct elements summaries are “mergeable”. No need to share lists of distinct elements if those elements are stored on different machines. Just share minimum hash value.

**Implementations:** Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

**Use Case:** Exploratory SQL-like queries on tables with 100's of billions of rows.

- **Count** number of **distinct** users in Germany that made at least one search containing the word 'auto' in the last month.
- **Count** number of **distinct** subject lines in emails sent by users that have registered in the last week.

## HYPERLOGLOG IN PRACTICE

**Implementations:** Google PowerDrill, Facebook Presto, Twitter Algebird, Amazon Redshift.

**Use Case:** Exploratory SQL-like queries on tables with 100's of billions of rows.

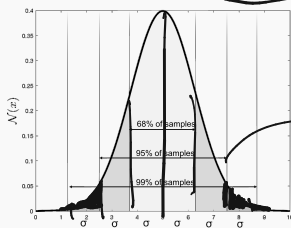
- **Count** number of **distinct** users in Germany that made at least one search containing the word 'auto' in the last month.
- **Count** number of **distinct** subject lines in emails sent by users that have registered in the last week.

Answering a query requires a (distributed) linear scan over the database: 2 seconds in Google's distributed implementation.

**Google Paper:** "Processing a Trillion Cells per Mouse Click"

**Motivating question:** Is Chebyshev's Inequality tight?

It is the worst case, but often not in reality.



$$\Pr(|X - \mathbb{E}(X)| > \underline{k\sigma}) \leq \frac{1}{k^2}$$

$$e^{-x^2/2\sigma^2}$$

68-95-99 rule for Gaussian bell-curve.  $X \sim N(\underline{0}, \underline{\sigma^2})$

**Chebyshev's Inequality:**

$$\Pr(|X - \mathbb{E}[X]| \geq 1\sigma) \leq \underline{100\%}$$

$$\Pr(|X - \mathbb{E}[X]| \geq 2\sigma) \leq \underline{25\%} \quad \frac{1}{2^2}$$

$$\Pr(|X - \mathbb{E}[X]| \geq 3\sigma) \leq \underline{11\%}$$

$$\Pr(|X - \mathbb{E}[X]| \geq 4\sigma) \leq \underline{6\%}$$

**Truth:**

$$\Pr(|X - \mathbb{E}[X]| \geq \underline{1\sigma}) \approx 32\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 2\sigma) \approx 5\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 3\sigma) \approx 1\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 4\sigma) \approx .01\%$$

## GAUSSIAN CONCENTRATION

$X \sim \mathcal{N}(\mu, \sigma^2)$  has probability density function (PDF)  $p$  with:

$$\underline{p(\mu \pm x)} = \frac{1}{\sigma\sqrt{2\pi}} \underline{e^{-x^2/2\sigma^2}}$$

### Lemma (Gaussian Tail Bound)

For  $\underline{X} \sim \mathcal{N}(\underline{\mu}, \underline{\sigma^2})$ :

$$\Pr[\underline{|X - \mathbb{E}X|} \geq \underline{k \cdot \sigma}] \leq \underline{2e^{-k^2/2}}.$$

Compare this to:

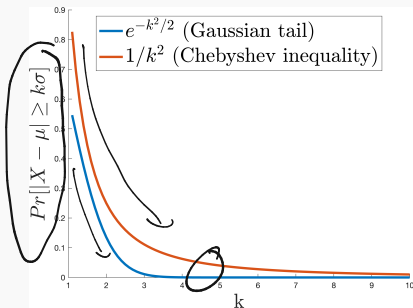
### Lemma (Chebyshev's Inequality)

For  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

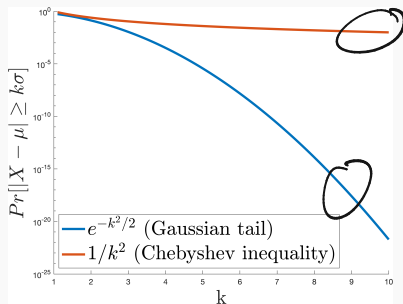
$$\Pr[|X - \mathbb{E}X| \geq k \cdot \sigma] \leq \frac{1}{k^2}$$



# GAUSSIAN CONCENTRATION



Standard y-scale.



Logarithmic y-scale.

**Takeaway:** Gaussian random variables concentrate much tighter around their expectation than variance alone predicts (i.e., than Chebyshev's inequality predicts).

Why does this matter for algorithm design?

# CENTRAL LIMIT THEOREM

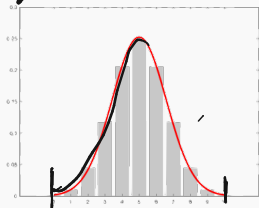
## Theorem (CLT – Informal)

Any sum of mutually independent, (identically distributed) r.v.'s  $\underline{X_1}, \dots, \underline{X_n}$  with mean  $\underline{\mu}$  and finite variance  $\underline{\sigma^2}$  converges to a Gaussian r.v. with mean  $n \cdot \mu$  and variance  $n \cdot \sigma^2$ , as  $n \rightarrow \infty$ .

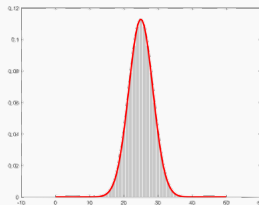
$$S = \sum_{i=1}^n X_i \implies \mathcal{N}(n \cdot \mu, n \cdot \sigma^2).$$

$$E\left(\sum_{i=1}^n X_i\right) = n \cdot \mu$$

$$\text{Var}(\sum X_i) = n \cdot \sigma^2$$



(a) Distribution of # of heads after 10 coin flips, compared to a Gaussian.



(b) Distribution of # of heads after 50 coin flips, compared to a Gaussian.

Recall:

## Definition (Mutual Independence)

Random variables  $X_1, \dots, X_n$  are mutually independent if, for all possible values  $v_1, \dots, v_n$ ,

$$\Pr[X_1 = v_1, \dots, X_n = v_n] = \Pr[X_1 = v_1] \cdot \dots \cdot \Pr[X_n = v_n]$$

Strictly stronger than pairwise independence.

## EXERCISE

If I flip a fair coin 100 times, lower bound the chance I get between 30 and 70 heads?

Let's approximate the probability by assuming the limit of the CLT holds exactly – i.e., that this sum looks exactly like a Gaussian random variable.

Lemma (Gaussian Tail Bound)

For  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$\Pr(|S - \mathbb{E}[S]| \geq 20) \leq 2e^{-4^2/2} =$$

$$\Pr[|X - \mathbb{E}X| \geq k \cdot \sigma] \leq 2e^{-k^2/2}.$$

$$S = \sum_{i=1}^n \mathbb{1}(\text{ith coin is heads}) \quad \mathbb{E}[S] = 50 \quad \text{Var}[S] = 25$$

$\downarrow$   
# of heads

$2e^{-8} = .06\%$  Chebyshev's inequality gave a bound of  $6.25\%$ .

These back-of-the-envelop calculations can be made rigorous! Lots of different “versions” of bound which do so.

- Chernoff bound   ?
- Bernstein bound   )
- Hoeffding bound   )
- ...

Different assumptions on random variables (e.g. binary vs. bounded), different forms (additive vs. multiplicative error), etc. Wikipedia is your friend.

## QUANTITATIVE VERSIONS OF THE CLT

### Theorem (Chernoff Bound)

Let  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  be independent  $\{0, 1\}$ -valued random variables and let  $p_i = \mathbb{E}[X_i]$ , where  $0 < p_i < 1$ . Then the sum  $\underline{S} = \sum_{i=1}^n X_i$ , which has mean  $\mu = \sum_{i=1}^n p_i$ , satisfies

$$\Pr[\underline{S} \geq (1 + \epsilon)\mu] \leq e^{\frac{-\epsilon^2 \mu}{2 + \epsilon}}$$

→ for  $\epsilon > 0$

and for  $0 < \epsilon < 1$

$$\Pr[S \leq (1 - \epsilon)\mu] \leq e^{\frac{-\epsilon^2 \mu}{2}}$$

$$\mathbb{E}[S] = \sum_{i=1}^n \mathbb{E}[X_i] = \sum_{i=1}^n p_i = \mu$$

# CHERNOFF BOUND

## Theorem (Chernoff Bound Corollary)

Let  $X_1, X_2, \dots, X_n$  be independent  $\{0, 1\}$ -valued random variables and let  $p_i = \mathbb{E}[X_i]$ , where  $0 < p_i < 1$ . Let  $S = \sum_{i=1}^n X_i$  and  $\mathbb{E}[S] = \mu$ . For  $\epsilon \in (0, 1)$ ,

$$\Pr[|S - \mu| \geq \epsilon \mu] \leq \frac{2e^{-\epsilon^2 \mu / 3}}{2e^{-k^2/3}}$$

$\nearrow k = 6$   
 $\epsilon \mu = 6$

$$\epsilon \mu = 6$$

$$\epsilon \mu \approx \sqrt{\mu} \cdot k$$

$$\epsilon \approx \frac{k}{\sqrt{\mu}}$$

Why does this look like the Gaussian tail bound of  $\Pr[|S - \mu| \geq k \cdot \sigma] \lesssim 2e^{-k^2/2}$ ? What is  $\sigma(S)$ ?

$$\text{Var}(S) = \sum_{i=1}^n \text{Var}(X_i) = \sum_{i=1}^n \underbrace{\mathbb{E}[X_i^2]}_{p_i} - \underbrace{\mathbb{E}[X_i]^2}_{p_i^2} = \sum_{i=1}^n p_i - p_i^2 = O\left(\sum_{i=1}^n p_i\right) = O(\mu)$$

$\sigma \approx \sqrt{\mu}$   
 $\epsilon \approx \frac{6}{\sqrt{\mu}}$

## QUANTITATIVE VERSIONS OF THE CLT

### Theorem (Bernstein Inequality)

Let  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  be independent random variables with each  $X_i \in [-1, 1]$ . Let  $\underline{\mu}_i = \mathbb{E}[X_i]$  and  $\underline{\sigma}_i^2 = \text{Var}[X_i]$ . Let  $\mu = \sum_{i=1}^n \mu_i$  and  $\sigma^2 = \sum_{i=1}^n \sigma_i^2$ . Then, for  $k \leq \frac{1}{2}\sigma$ ,  $S = \sum_{i=1}^n X_i$  satisfies

$$\Pr[|S - \mu| > k \cdot \sigma] \leq 2e^{-k^2/4}.$$

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

$$\mathbb{E}[S] = \mu = \sum_{i=1}^n \mu_i$$

$$\text{Var}[S] = \sigma^2 = \sum_{i=1}^n \sigma_i^2$$



## QUANTITATIVE VERSIONS OF THE CLT

### Theorem (Hoeffding Inequality)

Let  $X_1, X_2, \dots, X_n$  be independent random variables with each  $X_i \in [a_i, b_i]$ . Let  $\mu_i = \mathbb{E}[X_i]$  and  $\mu = \sum_{i=1}^n \mu_i$ . Then, for any  $k > 0$ ,  $S = \sum_{i=1}^n X_i$  satisfies:

$$k = k' \sqrt{\sum_{i=1}^n (b_i - a_i)^2}$$

$$\Pr[|S - \mu| > k] \leq 2e^{\frac{-2k^2}{\sum_{i=1}^n (b_i - a_i)^2}}.$$



$$\Pr[|S - \mu| \geq k' \sqrt{\sum_{i=1}^n (b_i - a_i)^2}] \leq \underline{2e^{-2k'^2}}$$

$$\mathbb{E}[|S - \mathbb{E}X|^2] = \sigma^2$$

$$\Pr[|S - \mathbb{E}X|^2 \geq \alpha \sigma^2] \leq \frac{1}{\alpha}$$

$$\Pr[|S - \mathbb{E}X| \geq \sqrt{\alpha} \sigma] \leq \frac{1}{\alpha}$$

## HOW ARE THESE BOUNDS PROVEN?

Return at 3:45.

Variance is a natural measure of central tendency, but there are others.

$q^{\text{th}}$  central moment:  $\mathbb{E}[(X - \mathbb{E}X)^q]$

$q = 2$  gives the variance. Proof of Chebyshev's applies Markov's inequality to the random variable  $(X - \mathbb{E}X)^2$ .

**Idea in brief:** Apply Markov's inequality to  $\mathbb{E}[(X - \mathbb{E}X)^q]$  for larger  $q$ , or more generally to  $f(X - \mathbb{E}X)$  for some other non-negative function  $f$ . E.g., to  $\exp(X - \mathbb{E}X)$ . Doing so requires higher-order independence.

$$\mathbb{E}[X_1, X_2, X_3, X_4] = \mathbb{E}[X_1] \mathbb{E}[X_2] \mathbb{E}[X_3] \mathbb{E}[X_4]$$

## EXERCISE

If I flip a fair coin 100 times, lower bound the chance I get between 30 and 70 heads?  $\rightarrow$  # of heads

Corollary of Chernoff bound: Let  $S = \sum_{i=1}^n X_i$  and  $\mu = \mathbb{E}[S]$ . For  $0 < \epsilon < 1$ ,

$$\Pr[|S - \mu| \geq \epsilon\mu] \leq \underline{\underline{2e^{-\epsilon^2\mu/3}}}$$

Here  $X_i = \mathbb{1}[i^{\text{th}} \text{ flip is heads}]$ .

$$\Pr[|S - \mu| \geq 20] = \Pr[|S - \mu| \geq \frac{2}{5} \cdot \mu] \leq 2e^{-\frac{(2/5)^2 \cdot 50}{3}} = 0.139... \approx \underline{\underline{1.4\%}}$$

$$6.25\%$$

$$1.4\%$$

$$0.06\%$$

# CHERNOFF BOUND APPLICATION

→ 1 with probability  $b$ .

General Statement: Flip biased coin  $n$  times. i.e. the coin is heads with probability  $b$ . As long as  $n \geq O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$

Flip  $n$  times  
 $S = \sum_{i=1}^n X_i \rightarrow \{0, 1\}$

$$\Pr[|S - \underbrace{b \cdot n}_{E[S]}| \geq \epsilon n] \leq \delta$$

$$\sum_n \Pr\{|S - E[S]| \geq \alpha E[S]\} \leq 2e^{-\alpha^2 E[S] / 3}$$

Chernoff

$$\alpha b n = \epsilon n \quad \alpha = \frac{\epsilon}{b}$$

$$\begin{aligned} \Pr\{|S - E[S]| \geq \epsilon n\} &\leq 2e^{-\frac{\epsilon^2}{b^2} n / 3} \\ &\leq 2e^{-\epsilon^2 n / 3} = \delta \end{aligned}$$

Pay very little for higher probability – if you increase the number of coin flips by  $4x$ ,  $\delta$  goes from  $e^{-\epsilon^2 n / 3} = \delta/2$

$$1/10 \rightarrow 1/100 \rightarrow 1/10000$$

$$e^{-\epsilon^2 n / 3} = \log(\delta/2)$$

$$\epsilon^2 n / 3 = \log(2/\delta) \quad n = \frac{3 \log(2/\delta)}{\epsilon^2}$$

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

### Load balancing problem:

Suppose Google answers map search queries using servers  $A_1, \dots, A_q$ . Given a query like “new york to rhode island”, common practice is to choose a random hash function  $h \rightarrow \{1, \dots, q\}$  and to route this query to server:

$$A_h(\text{“new york to rhode island”})$$

**Goal:** Ensure that requests are distributed evenly, so no one server gets loaded with too many requests. We want to avoid downtime and slow responses to clients.

Why use a hash function instead of just distributing requests randomly?

# LOAD BALANCING

$$m/n$$

$n = \# \text{ of servers}$

Suppose we have  $n$  servers and  $m$  requests,  $x_1, \dots, x_m$ . Let  $s_i$  be the number of requests sent to server  $i \in \{1, \dots, n\}$ :

$\# \text{ of jobs sent to server } i$   $\rightarrow$   $s_i = \sum_{j=1}^m \mathbb{I}[h(x_j) = i]$   $\rightarrow$   $s_1, \dots, s_n$

Formally, our goal is to understand the value of maximum load on any server, which can be written as the random variable:

$$S = \max_{i \in \{1, \dots, n\}} s_i$$

## LOAD BALANCING

A good first step is to first think about expectations. If we have  $n$  servers and  $m$  requests, for any  $i \in \{1, \dots, n\}$ :

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

But it's unclear what the expectation of  $S = \max_{i \in \{1, \dots, n\}} S_i$  is...  
in particular,  $\mathbb{E}[S] \neq \max_{i \in \{1, \dots, n\}} \mathbb{E}[S_i]$ .

**Exercise:** Convince yourself that for two random variables  $A$  and  $B$ ,  $\mathbb{E}[\max(A, B)] \neq \max(\mathbb{E}[A], \mathbb{E}[B])$  even if those random variable are independent.

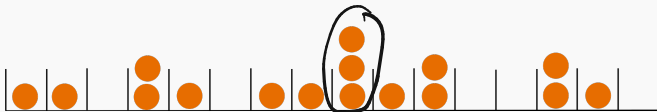
0 0  
1 1

coin 1 : 0 1 0 1  
coin 2 : 0 0 1 1

## SIMPLIFYING ASSUMPTIONS

**Number of servers:** To reduce notation and keep the math simple, let's assume that  $m = n$ . I.e., we have exactly the same number of servers and requests.

**Hash function:** Continue to assume a fully (uniformly) random hash function  $h$ .



Often called the “balls-into-bins” model.

$\mathbb{E}[s_i]$  = expected number of balls per bin  $= \frac{m}{n} = 1$ . We would like to prove a bound of the form:

$$\Pr[\max_i s_i \geq C] \leq \frac{1}{10}.$$

for as tight a value of  $C$ . I.e., something much better than  $C = n$ .



## BOUNDING A UNION OF EVENTS

**Goal:** Prove that for some  $C$ ,

$$\Pr[\max_i s_i \geq C] \leq \frac{1}{10}.$$

**Equivalent statement:** Prove that for some  $C$ ,

$$\Pr[(s_1 \geq C) \cup (s_2 \geq C) \cup \dots \cup (s_n \geq C)] \leq \frac{1}{10}.$$

These events are not independent, but we can apply union bound!

$$\leq \sum_{i=1}^n \Pr[s_i \geq C] \leq \sum_{i=1}^n \frac{1}{10n} = \frac{1}{10}$$

$n$  = number of balls and number of bins.  $s_i$  is number of balls in bin  $i$ .  $C$  = upper bound on maximum number of balls in any bin.

$$C_1 : 0 \quad 1 \quad 0 \quad 1$$

$$C_2 : 0 \quad 0 \quad 1 \quad 1$$

$$\Pr(C_1 \text{ or } C_2) \geq 1/2 \\ = 3/4$$

$$\Pr(C_1 \geq 1/2) \\ + \Pr(C_2 \geq 1/2) \\ = 1$$

## APPLICATION OF UNION BOUND

We want to prove that:

$$\Pr[\max_i s_i \geq C] = \Pr[(s_1 \geq C) \cup (s_2 \geq C) \cup \dots \cup (s_n \geq C)] \leq \frac{1}{10}.$$

To do so, it suffices to prove that for all  $i$ :

$$\Pr[s_i \geq C] \leq \frac{1}{10n}.$$

Why? Because then by the union bound,

$$\begin{aligned} \Pr[\max_i s_i \geq C] &\leq \sum_{i=1}^n \Pr[s_i \geq C] \quad (\text{Union bound}) \\ &\leq \sum_{i=1}^n \frac{1}{10n} = \frac{1}{10}. \quad \square \end{aligned}$$

$n$  = number of balls and number of bins.  $s_i$  is number of balls in bin  $i$ .

## NEW GOAL

$$\Pr[s_i \geq k \mid \mathbb{E}(s_i)] \leq 1/k$$

$$k = 10n$$

$$\Pr[s_i \geq k] \leq 1/k$$

$$\Pr[s_i \geq \underline{10n}] \leq \frac{1}{10n}$$

Prove that for some  $C$ ,

$$\Pr[\underline{s_i} \geq C] \leq \left( \frac{1}{10n} \right)$$

Let's try doing this with Markov's, Chebyshev, and exponential concentration.

## ATTEMPT WITH MARKOV'S INEQUALITY

**Goal:** Prove that  $\Pr[s_i \geq C] \leq \frac{1}{10n}$ .

- **Step 1.** Verify we can apply Markov's:  $s_i$  takes on non-negative values only. Good to go!
- **Step 2.** Apply Markov's:  $\Pr[s_i \geq C] \leq \frac{\mathbb{E}[s_i]}{C} = \frac{1}{C}$ .

To prove our target statement, need to see  $C = 10n$ .

Meaningless! There are only  $n$  balls, so of course there can't be more than  $10n$  in the most overloaded bin.

$n$  = number of balls and number of bins.  $s_i$  is number of balls in bin  $i$ .  $\mathbb{E}[s_i] = 1$ .  $C$  = upper bound on maximum number of balls in any bin. **Markov's inequality:** for positive r.v.  $X$ ,  $\Pr[X \geq t] \leq \mathbb{E}[X]/t$ .

## ATTEMPT WITH CHEBYSHEV'S INEQUALITY

**Goal:** Prove that  $\Pr[s_i \geq C] \leq \frac{1}{10n}$ .

- **Step 1.** To apply Chebyshev's inequality, we need to understand  $\sigma^2 = \text{Var}[s_i]$ .

Use linearity of variance. Let  $s_{i,j}$  be a  $\{0, 1\}$  indicator random variable for the event that ball  $j$  falls in bin  $i$ . We have:

*# of balls in bin i*  $\rightarrow$   $s_i = \sum_{j=1}^n s_{i,j}$   $\leftarrow$   *$\mathbb{I}(\text{ball } j \in \text{bin } i)$*

$n$  = number of balls and number of bins.  $s_i$  is number of balls in bin  $i$ .  $\mathbb{E}[s_i] = 1$ .  $C$  = upper bound on max number of balls in bin.

$$\text{Var}[s_i] = \sum_{j=1}^n \text{Var}[s_{i,j}] = \sum_{j=1}^n \frac{1}{n} - \frac{1}{n^2} \leq \sum_{j=1}^n \frac{1}{n} = \textcircled{1}$$

## VARIANCE ANALYSIS

$$s_{i,j} = \begin{cases} 1 & \text{with probability } \frac{1}{n} \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathbb{E}[s_{i,j}] =$$

$$\mathbb{E}[s_{i,j}^2] =$$

So:

$$\text{Var}[s_i] = \text{Var} \left[ \sum_{j=1}^n s_{i,j} \right] =$$

$n$  = number of balls and number of bins.  $s_{i,j}$  is event ball  $j$  lands in bin  $i$ .

## APPLYING CHEBYSHEV'S

Goal: Prove that  $\Pr[s_i \geq C] \leq \frac{1}{10n}$ .

$$6^2 = 1 \quad 6 = 1$$

Step 1. To apply Chebyshev's inequality, we need to understand  $\sigma^2 = \text{Var}[s_i] = 1$

$$\text{Var}[s_i] = \sum_{j=1}^n \text{Var}[s_{i,j}] = \sum_{j=1}^n \frac{1}{n} - \frac{1}{n^2} = 1 - \frac{1}{n} \leq 1.$$

Step 2. Apply Chebyshev's inequality:

$$\Pr[|s_i - \mathbb{E}[s_i]| \geq k \cdot 1] \leq \frac{1}{k^2}$$

$$\Pr[s_i \geq \sqrt{10n} + 1] \leq \frac{1}{10n}.$$

$n$  = number of balls and number of bins.  $s_i$  = number of balls in bin  $i$ .  $s_{i,j}$  is event ball  $j$  lands in bin  $i$ .  $\mathbb{E}[s_i] = 1$ .

$$k = \sqrt{10n} \quad \Pr[|s_i - 1| \geq \sqrt{10n}] \leq \frac{1}{10n}$$

## APPLYING CHEBYSHEV'S

**Goal:** Prove that  $\Pr[s_i \geq C] \leq \frac{1}{10n}$ .

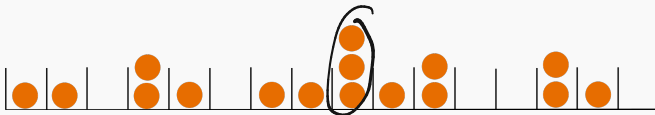
We just proved that, for any  $k$ :  $\Pr[|s_i - 1| \geq k] \leq \frac{1}{k^2}$ .

$n$  = number of balls and number of bins.  $s_i$  is number of balls in bin  $i$ .  $C$  = upper bound on maximum number of balls in any bin.



## FINAL RESULT FOR CHEBYSHEV'S

When hashing  $n$  balls into  $n$  bins, the maximum bin contains  $\Theta(\sqrt{n})$  balls with probability  $\frac{9}{10}$ .



Much better than the trivial bound of  $n$ !

## ATTEMPT WITH EXPONENTIAL CONCENTRATION

**Goal:** Prove that  $\Pr[s_i \geq C] \leq \frac{1}{10n}$ .

**Recall:**  $s_i = \sum_{j=1}^n s_{i,j}$ , where  $s_{i,j} = \mathbb{1}[\text{ball } j \text{ lands in bin } i]$ .



What bound might we use?

*Chernoff*

$\mathbb{1}[\text{ball } j \text{ lands in bin } i]$

## ATTEMPT WITH EXPONENTIAL CONCENTRATION

### Theorem (Chernoff Bound)

Let  $X_1, X_2, \dots, X_n$  be independent  $\{0, 1\}$ -valued random variables and let  $p_i = \mathbb{E}[X_i]$ , where  $0 < p_i < 1$ . Then the sum  $S = \sum_{j=1}^n X_j$ , which has mean  $\mu = \sum_{j=1}^n p_j$ , satisfies

$$\Pr[S \geq (1 + \epsilon)\mu] \leq e^{\frac{-\epsilon^2 \mu}{2 + \epsilon}} = \frac{1}{10n}$$

Apply with  $S = \sum_{j=1}^n X_j$

$$\Pr[S \geq (1 + c \log n)\mu] \leq e^{-\frac{c^2 \log^2 n}{2 + c \log n}} = e^{-\frac{c \log^2 n}{2/c + 1 + \log n}} \leq e^{-\frac{c \log^2 n}{2 \log n}} = e^{-\frac{c \log n}{2}} = \left(\frac{1}{n}\right)^{c/2} \leq \frac{1}{10n}$$

$\underbrace{c \log n}_{\alpha(\log n)}$

## LOAD BALANCING

$$1 - \frac{1}{n^2} \quad O(\log n) \quad \frac{1}{n^2} \quad \frac{1}{n^2} \cdot n = \frac{1}{n}$$

So max load for randomized load balancing is  $O(\log n)$ . Best we could prove with Chebyshev's was  $O(\sqrt{n})$ .

$$\{f_{\max}\} \leq O(\log n) + 1/n = O(\log n)$$

## POWER OF TWO CHOICES

$$(\log \log \log n)$$

Power of 2 Choices: Instead of assigning job to random server, choose 2 random servers and assign to the least loaded. With probability  $1/10$  the maximum load is bounded by:

- (a)  $O(\log n)$     (b)  $O(\sqrt{\log n})$     (c)  $O(\log \log n)$     (d)  $O(1)$

