

Tuesday, February 10

- Thank you for coming to GEMS/my talk!
- CCMS Applied Math seminar Mondays: 4:15 pm Estella 1021
- CCMS Colloquium Fridays: 11am Davidson Lecture Hall
- No OH Thursday, DM to meet Friday!

Plan

- Clean up discrete elements
- Load balancing

Distinct Elements

Wp $1-\delta$ in $O\left(\frac{\log D}{\epsilon^2 \delta}\right)$ space:

$$(1-\epsilon)\mu \leq \bar{S} \leq (1+\epsilon)\mu$$

$$\Leftrightarrow \frac{1}{(1+\epsilon)\mu} \leq \frac{1}{\bar{S}} \leq \frac{1}{(1-\epsilon)\mu}$$

$$1-2\epsilon \leq \frac{1}{1+\epsilon} ; \frac{1}{1-\epsilon} \leq 1+2\epsilon \text{ by Desmos}$$

$$\Leftrightarrow (1-2\epsilon) \frac{1}{\mu} - 1 \leq \frac{1}{\bar{S}} - 1 \leq (1+2\epsilon) \frac{1}{\mu} - 1$$

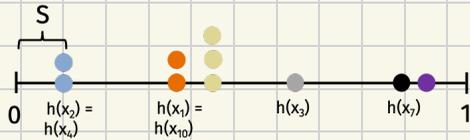
$$\mu = \frac{1}{2}$$

$$\Leftrightarrow (1-4\epsilon) \left(\frac{1}{\mu} - 1\right) \leq \hat{D} \leq (1+4\epsilon) \left(\frac{1}{\mu} - 1\right)$$

$$\Leftrightarrow (1-4\epsilon) D \leq \hat{D} \leq (1+4\epsilon) D$$

$$(1-\epsilon) D \leq \hat{D} \leq (1+\epsilon) D$$

where we hide 4 in big-O



$$\bar{S} \leftarrow \sum_{j=1}^K \frac{S_j}{\kappa}$$

$$D \leftarrow \frac{1}{\delta} - 1$$

Load Balancing

Goal: Distribute load evenly between servers using hash function

e.g., Google maps routing

Advantage of hash: cache redundant queries

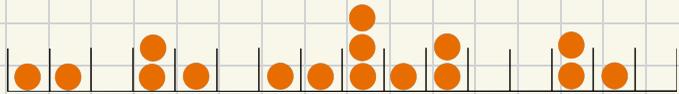
n servers and m (unique) queries

$$\begin{aligned} S_i &= \# \text{ queries sent to server } i \\ &= \sum_{j=1}^m \mathbb{1}[h(x_j) = i] \end{aligned}$$

Bound $S = \max_{i \in [n]} S_i$

- $\mathbb{E}[S_i] =$
- $\mathbb{E}[S] \stackrel{2.}{=} \max_{i \in [n]} \mathbb{E}[S_i]$

Let's assume $m=n$, so $\mathbb{E}[S_i]=1$



Goal: Prove

$$\Pr(\max_i s_i \geq C) \leq 1/10$$

$$\Leftrightarrow \Pr(s_1 \geq C \cup s_2 \geq C \cup \dots \cup s_n \geq C) \leq 1/10$$

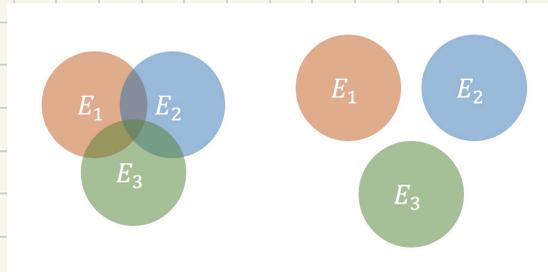
By union bound,

$$\Pr(s_1 \geq C \cup \dots \cup s_n \geq C)$$

$$\leq \sum_{i=1}^n \Pr(s_i \geq C) = n \Pr(s_i \geq C) \stackrel{\text{want}}{\leq} n \cdot \frac{1}{10n}$$

Union Bound

Proof by picture:



Proof by Markov's:

$$\begin{aligned} \Pr(E_1 \cup \dots \cup E_n) &= \Pr\left(\sum_{i=1}^n \mathbb{1}[E_i] \geq 1\right) \\ &\leq \mathbb{E} \sum_{i=1}^n \mathbb{1}[E_i] = \sum_{i=1}^n \Pr(E_i) \end{aligned}$$

Prove that

$$\Pr(s_i \geq c) \leq \frac{1}{10n}$$

- Markov's ?
- Chebyshev's ?
- What c ?

Thursday, February 12

Plan

Concentration Inequalities

Revisit coins + load balancing

Chebyshev's Inequality

Let X rv with $\mu = \mathbb{E}[X]$, $\sigma^2 = \text{Var}(X)$

For $k > 0$,

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

On load balancing, showed

max load $\geq \Theta(\sqrt{n})$ whp.

Using tighter inequality, we

can show max load $\geq O(\log n)$ whp.

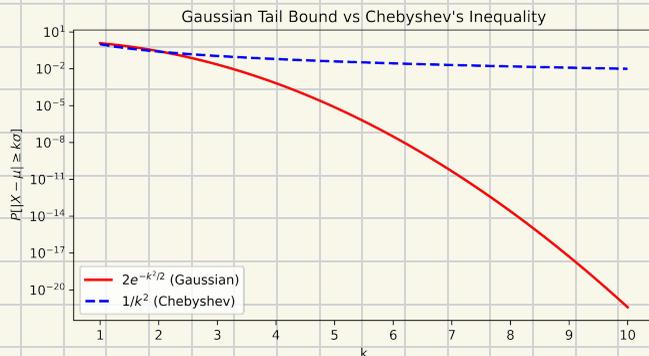
Q: When is Chebyshev loose?

Gaussian Tail Bound

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

For $k > 0$,

$$\Pr(|X - \mu| > k\sigma) \approx 2e^{-k^2/2}$$



Q: Is Chebyshev always loose?

Q: For which RV is Chebyshev's loose?

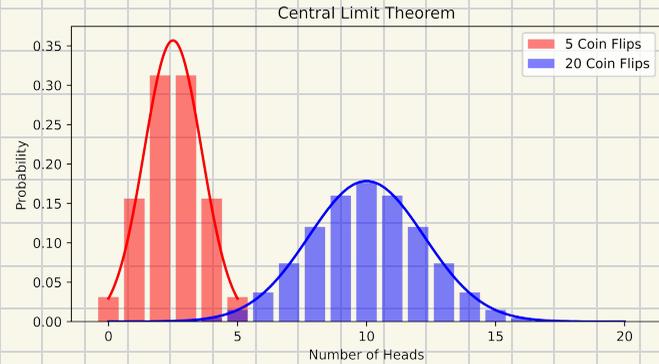
Central Limit Theorem:

The sum of mutually independent X_1, \dots, X_n with $\mu = E[X_i]$, $\sigma^2 = \text{Var}(X_i)$ converges to $\mathcal{N}(n\mu, n\sigma^2)$ i.e.,

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

↑ also holds for rv w/ different mean, var

↑ Today: When does CLT "kick in"?



$$n = 100 \quad \mu = 50 \quad \sigma^2 = 25 \quad k = 4$$

Chebyshev's

$$\Pr(|X - \mu| \geq k\sigma) \leq 1/16 = 0.0625$$

Exponential Tail Bound

$$\Pr(|X - \mu| \geq k\sigma) \leq 2e^{-k^2/2} \approx 0.006$$

Exponential Concentration Inequalities

X_1, \dots, X_n indep

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

↳ Make assumptions on X_i :
stronger assumption = stronger bound

↳ Proved using clever applications
of Markov's

Chernoff Bound

Binary $X_i \in \{0, 1\}$. For $0 < \epsilon < 1$,

$$\Pr(S - \mu \geq \epsilon \mu) \leq 2 \exp\left(-\frac{\epsilon^2 \mu}{3}\right)$$

For $\epsilon > 0$,

$$\Pr(S \geq (1 + \epsilon)\mu) \leq \exp\left(-\frac{\epsilon^2 \mu}{2 + \epsilon}\right)$$

Bernstein Inequality

$$X_i \in [-1, 1]. \quad \sigma^2 = \sum_{i=1}^n \text{Var}(X_i).$$

For $k \leq \sigma/2$,

$$\Pr(|S - \mu| > k\sigma) \leq 2 \exp\left(-\frac{k^2}{4}\right)$$

$$\begin{aligned} \text{Var}(X_i) &= \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] \\ &\leq (b_i - a_i)^2 \end{aligned}$$

Hoeffding's Inequality

$X_i \in [a_i, b_i]$

$$\Pr(|S - \mu| > k) \leq 2 \exp\left(-\frac{2k^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

$$\text{Var } X_i \leq \mathbb{E}[X_i]$$

Revisiting Coin Flips

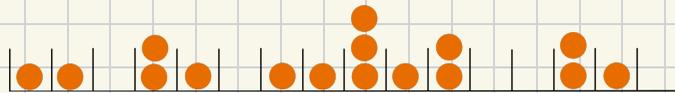
$$X_i \in \{0, 1\} \quad \mathbb{E}[X_i] = b \quad \forall i$$

Choose $\delta > 0$, $\epsilon > 0$, then $n \geq \frac{3 \log(2/\delta)}{\epsilon^2}$.

$$\Pr(|S - bn| \geq \epsilon n) \leq \delta$$

Proof:

Revisiting Load Balancing



$$S_i = \sum_{j=1}^n \mathbb{1}[j \text{ to } i]$$

$$S = \max_i S_i$$

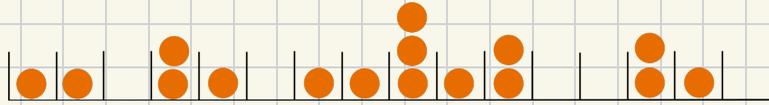
$$\Pr(S \geq c) \leq \frac{1}{10}$$

$$\Leftrightarrow \Pr(S_i \geq c) \leq 1/10n$$

with Chebyshev's, $c = O(\sqrt{n})$

Q: Can we do better?

Power of Two Choices



use two hash functions,
choose least occupied

Then $\Pr(S \geq c) \leq 1/10$ for $c = \log n$

$$c = \log \log n$$

$$c = \log \log \log n$$

How about power of three choices?