

Week 5: Chernoff Bounds and Sampling

February 12, 2019

Last week: Hashing

- (1) Chaining
- (2) Linear Probing
- (3) Cuckoo Hashing

Today

1) Chernoff-Hoeffding Bounds

2) Examples: Coin flipping

load balancing

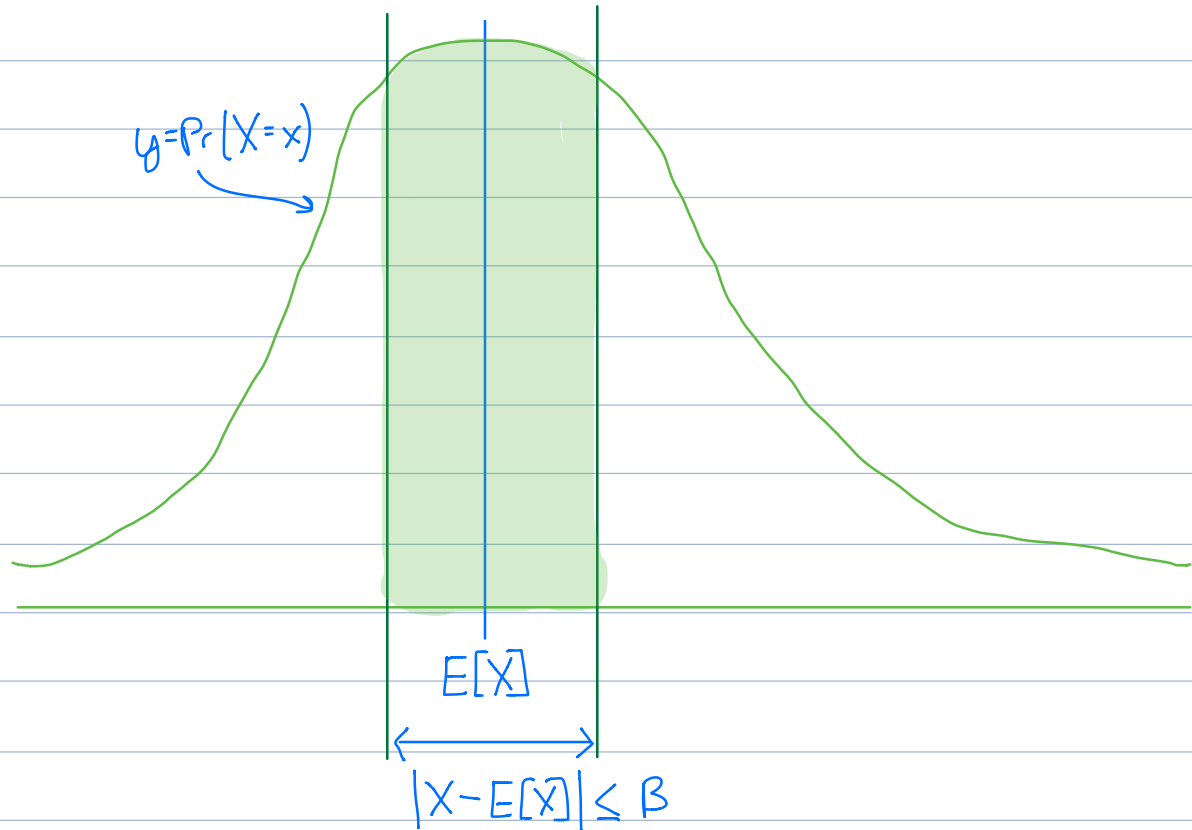
Sampling

approximate counting

Deviation from the Mean

X = random variable

$E[X]$ = expected value



Goal: Show $\Pr(|X - E[X]| > B)$ is small.

Ex. Flipping coins

→ flip n coins

→ $\Pr(\text{heads}) = 1/2$

→ $X = \# \text{ heads}$

1) $E[X] = \frac{n}{2} \leftarrow \mu$

→ $X_j = 1$ if j^{th} coin is heads
= 0 otherwise

→ $E[X_j] = \Pr(X_j = 1) = 1/2$

→ $E[X] = E[\sum X_j] = \sum E[X_j] = n \cdot \frac{1}{2}$

2) Markov's Inequality: $\Pr(X \geq \frac{3}{4}n) \leq \frac{n/2}{\frac{3}{4}n} = \frac{2}{3}$

→ not so useful here!

3) Chebyshev

→ $\text{Var}(X_j) = p(1-p) = \frac{1}{4}$

→ $\text{Var}(X) = \frac{n}{4}$

→ $\Pr(X \geq \frac{3}{4}n) \leq \Pr(|X - \mu| \geq \frac{n}{4}) \leq \frac{\text{Var}(X)}{(\frac{n}{4})^2}$

$\leq \frac{\frac{n}{4}}{(\frac{n}{4})^2} \leq \frac{4}{n}$

Better... What about $\Pr(|X - \mu| \geq \epsilon n)$??

Chernoff Bound: X_1, X_2, \dots, X_n independent
 $X_i \in \{0, 1\}$
 $\Pr(X_i) = p$
 $\mu = E[X] = np$

1) General Version:

$$\frac{\delta > 0}{\Pr(X \geq (1+\delta)\mu) \leq \left(\frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu}$$

$$\frac{0 < \delta < 1}{\Pr(X \leq (1-\delta)\mu) \leq \left(\frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} \right)^\mu}$$

$$\text{Ex: } \Pr(X \geq e\mu) \leq \left(\frac{e^{e-1}}{e^e} \right)^\mu \leq e^{-\mu}$$

2) Simplified versions

$$\delta > 0$$

$$\Pr(X \geq (1+\delta)M) \leq e^{-\left(\delta^2/2+\delta\right)M}$$

$$\Pr(X \leq (1-\delta)M) \leq e^{-M\delta^2/2} \leq e^{-M\delta^2/3}$$

$$\delta \leq 1$$

$$\Pr(X \leq (1-\delta)M) \leq e^{-M\delta^2/3}$$



These are what I memorize!

Coin flipping.

$$\begin{aligned} \Pr(X \geq \frac{7}{4}n) &\leq \Pr(X \geq (1+\frac{1}{2})\frac{n}{2}) \\ &\leq \Pr(X \geq (1+\frac{1}{2})M) \leq e^{-\frac{n}{2} \cdot (\frac{1}{2})^2 \cdot \frac{1}{3}} \\ &\leq \frac{1}{e^{-n/24}} \end{aligned}$$

Decreases exponentially in n !! Very unlikely...

Tighter?

$$\Pr(X \geq \mu + \lambda) = \Pr(X \geq (1 + \frac{\lambda}{\mu})\mu)$$

$$\leq e^{-\frac{\mu}{3} \cdot (\frac{\lambda}{\mu})^2}$$

$\lambda \leq \mu$:

$$\leq e^{-\frac{\lambda^2}{3\mu}}$$

$$\text{Choose } \lambda = \sqrt{3C\mu \ln(n)} = \sqrt{\frac{3}{2} C n \ln(n)}$$

$$\leq e^{-\frac{3C\mu \ln(n)}{3\mu}}$$

$$\leq e^{-C \ln(n)}$$

$$\leq \left(\frac{1}{n}\right)^C$$

$$\text{Same idea: } \Pr(X \leq \mu - \lambda) \leq \left(\frac{1}{n}\right)^C$$

Conclusion: With high probability (w.h.p.)
of heads is $\frac{n}{2} \pm \Theta(\sqrt{n \log n})$

$$\begin{aligned} \Pr(\text{fail}) &\leq \Pr(X \geq \mu + \lambda) + \Pr(X \leq \mu - \lambda) \\ &\leq \frac{1}{n^c} + \frac{1}{n^c} = \frac{2}{n^c} \end{aligned}$$

Proof of Chernoff Bound

$$\textcircled{1} \Pr[X \geq (1+\delta)M] = \Pr[tX \geq t(1+\delta)M]$$

Assume $t > 0$

$$= \Pr[e^{tX} \geq e^{t(1+\delta)M}]$$

$\textcircled{2}$ Markov's Inequality

$$\Pr[e^{tX} \geq e^{t(1+\delta)M}] \leq \frac{E[e^{tX}]}{e^{t(1+\delta)M}}$$

$\textcircled{3}$ Independence

$$E[e^{tX}] = E[e^{t \sum x_j}] = E\left[\prod_{j=1}^n e^{tx_j}\right]$$

$$\text{independence} \rightarrow = \prod_{j=1}^n E[e^{tx_j}]$$

④ Approximation

$$E[e^{tx_j}] = pe^t + (1-p)e^0$$

$$= 1 - p(1 - e^t) \quad x = p(1 - e^t)$$

$$\leq e^{-p(1 - e^t)} \quad e^{-x} \geq 1 - x$$

⑤ Conclusion:

$$\prod_{j=1}^n E[e^{tx_j}] \leq (e^{-p(1 - e^t)})^n \leq e^{-pn(1 - e^t)} \\ \leq e^{-\mu(1 - e^t)}$$

$$\text{Choose } t = \ln(1 + \delta) \leq e^{\mu\delta}$$

$$\Pr(X \geq (1 + \delta)\mu) \leq \frac{e^{\mu\delta}}{e^{t(1 + \delta)\mu}} \leq \frac{e^{\mu\delta}}{(1 + \delta)^{(1 + \delta)\mu}}$$

$$\leq \left(\frac{e^\delta}{(1 + \delta)^{(1 + \delta)}} \right)^\mu$$

Variations:

① X_j not identical: $\Pr(X_j=1) = p_j$
 $\mu = E[X] = \sum p_j$

$$\prod_{j=1}^n E[e^{tX_j}] \leq \left(e^{-p(1-e^t)} \right)^n \leq e^{-pn(1-e^t)}$$

$$\uparrow \leq \prod_{j=1}^n e^{-p_j(1-e^t)}$$

$$= e^{-\sum p_j(1-e^t)}$$

$$= e^{-\mu(1-e^t)}$$

OK ✓

② X_j not indicator: $X_j \in [0,1]$ ← bounded!
e.g. $X_j = 0.75$

$$E[e^{tX_j}] = pe^t + (1-p)e^0$$

Hoeffding's Lemma shows:

$$E[e^{tX}] \leq e^{tE[X]} e^{t^2/8}$$

So it still works...

③ X_j not in $[0, 1]$: $X_j \in [0, B]$ $\leftarrow B \geq 1$

Define $y_j = \frac{x_j}{B}$, $Y = \sum y_j = \frac{1}{B} \sum X_j$

$$\mu_y = E[y_j] = \frac{1}{B} E[X_j] = \frac{1}{B} \mu_x$$

$$\Pr(\sum X_j \geq (1+\delta)\mu_x) \leq \Pr(B \sum y_j \geq (1+\delta) B \mu_y)$$

$$\leq \Pr(\sum y_j \geq (1+\delta) \mu_y)$$

$$\leq e^{-\mu_y \delta^2 / 3}$$

$$\leq e^{-\frac{\mu_x \delta^2}{3B^2}} \leftarrow \text{reduces exponent by } B^2$$

OK if B is constant or small compared to μ_x

④ Two-sided error:

$$\Pr[|X - \mu| \geq \delta \mu] \leq \Pr(X \geq (1+\delta)\mu)$$

+ \leftarrow union bound

$$\Pr(X \leq (1-\delta)\mu)$$

$$\leq 2e^{-\mu \delta^2 / 3}$$

④ Absolute bound:

$$\Pr[|X - \mu| \geq \delta] \leq 2e^{-\frac{2\delta^2}{n}}$$

← pay factor of n
OK if $\delta^2 > n$

if $B > 1$: $\Pr[|X - \mu| \geq \delta] \leq 2e^{-\frac{2\delta^2}{nB^2}}$

Hoeffding's Bound

Ex: Coin flipping
 $\delta = \Theta(\sqrt{n \log n})$

⑤ Negative Correlation / Negative Association

$$\text{if } E[\pi e^{tX_i}] \leq \pi E[e^{tX_i}]$$

then... ok!

Works for balls in bins...

Ex Polling/Sampling

What fraction of Singaporeans like basketball?

Assume population size n , $p n$ like basketball
 $0 \leq p \leq 1$

Goal: find E such that:

$$\Pr[|E - pn| \leq \delta] \geq 1 - \epsilon$$

"Confidence interval $\pm \delta$ w.p. $\geq 1 - \epsilon$ "

"if $\delta = 0.02$ and $\epsilon = 0.1$: $\pm 2\%$ with 90% confidence"

$$[p - \delta, p + \delta] = \text{confidence interval}$$

Algo: Choose sample of S people.

Let $X_j = 1$ if person j likes basketball.

Return:

$$E = \left(\frac{1}{S}\right) \sum X_j$$

Analysis

$$\Pr[|E - p| \geq \delta] = \Pr[|\sum X_j - ps| \geq \delta s]$$

$$X = \sum X_j, \quad \mu = E[X] = ps$$

$$= \Pr[|X - \mu| \geq \left(\frac{\delta}{p}\right) \mu]$$

Chernoff: ① $\leq 2e^{-\frac{(\delta/p)^2 \cdot \mu}{2 + \delta/p}}$

$$\text{if } \delta < p$$

$$\leq 2e^{-\frac{\delta^2 s}{3p}}$$

$$\text{choose } s \geq \frac{3p}{\delta^2} \ln\left(\frac{2}{\epsilon}\right)$$

$$\leq \epsilon$$

need to know p ??

$$p \leq 1$$

$$\Rightarrow s \geq \frac{3}{\delta^2} \ln\left(\frac{2}{\epsilon}\right)$$

Sample Size does not depend on n !!

→ good for Singapore OR China!!

Hoeffding: $\mathbb{P}(\bar{Z}) \leq 2e^{\frac{-2\delta^2 s^2}{s}}$
 $\leq 2e^{-2\delta^2 s}$

Choose $s \geq \frac{1}{2\delta^2} \cdot \ln\left(\frac{2}{\epsilon}\right)$

$\leq \epsilon$

(OK if $\delta > p$?)

Load balancing: n jobs
 k servers

Algo: assign jobs randomly to servers

[Already analyzed $n=k$]

Fix a server. $X_j = 1$ if job j is on server.

$$E[X_j] = \Pr[X_j] = \frac{1}{k} \quad E[X] = \frac{n}{k}$$

3 cases

$$n \gg k$$

$$n > \Omega(k \log n)$$

$$n = \Theta(k)$$

$$\textcircled{1} \Pr[X \geq \frac{n}{k} + \lambda] \leq 2e^{-\frac{2\lambda^2}{n}}$$
$$\lambda \leq \sqrt{\frac{cn \log n}{2}}$$

$$\begin{aligned} \text{or: } \Pr[X \geq (1 + \frac{\lambda}{\frac{n}{k}}) \frac{n}{k}] \\ \leq e^{-\frac{\lambda^2 k^2}{n} \cdot \frac{n}{k} \cdot \frac{1}{2}} \\ \leq e^{-\frac{\lambda^2 k}{2n}} \\ \lambda = \sqrt{\frac{cn \log n}{2}} \end{aligned}$$

$$\leq \frac{2}{n^c}$$

With Prob $\geq 1 - \frac{2}{n^c}$, no server has more than

$$\frac{n}{k} + \sqrt{\frac{c}{2} n \log n} \text{ tasks}$$

$$\textcircled{2} \Pr[X \geq e \frac{n}{K}] \leq e^{-\frac{n}{K}}$$

if $n > c K \log n$

$$\leq e^{-c \log n} \leq \frac{1}{n^c}$$

With Prob $\geq 1 - \frac{1}{n^{c-1}}$, no server has more than $\frac{en}{K}$ jobs

$$\begin{aligned} \textcircled{3} \Pr(X \geq \log(n) \frac{n}{K}) &\leq e^{-\frac{(\log n - 1)^2 \cdot \frac{n}{K}}{1 + \log n}} \\ &\leq e^{-\frac{\frac{1}{2} \log^2 n \cdot \frac{n}{K}}{2 \log n}} \\ &\leq e^{-\frac{n \log n}{4K}} \\ &\leq \left(\frac{1}{n}\right)^{-\frac{n}{4K}} \end{aligned}$$

If $n > 4cK$, no server has more than $\log(n)$ jobs w.p. $\geq 1 - \frac{1}{n^{c-1}}$

Coloring a Bipartite Graph

$$G = (U, V, E), \text{ bipartite}$$

Color each node in V red or blue to

$$\text{minimize: } \forall u \in U, |R(u) - B(u)|$$

\uparrow red neighbors \uparrow blue neighbors

or:

$$\geq \left(\frac{1}{2} - \delta\right) \text{ nbrs of } u \text{ blue}$$
$$\geq \left(\frac{1}{2} - \delta\right) \text{ nbrs of } u \text{ red}$$

Fix u : let $K = \deg(u)$

$Y_j = 1$ if j^{th} neighbor is blue

$R(u) = K - B(u)$

$$|R(u) - B(u)| = 2 \left| \frac{K}{2} - \sum Y_j \right|$$

$$\Pr \left(\left| \sum Y_j - \frac{K}{2} \right| \geq \lambda \right) \leq 2e^{-\frac{2\lambda^2}{K}}$$

$$\lambda = \sqrt{\frac{c}{2} K \ln(n)}$$

$$\leq 2 \cdot \frac{1}{n^c}$$

$$\Rightarrow (\text{Union bound over } U) \Pr(\text{any } u \text{ fails}) \leq \frac{2}{n^{c-1}} \leq \frac{1}{n^{c-2}}$$

$$E_{\delta}, \text{ if } \frac{K}{4} \geq \sqrt{\frac{c}{2} K \ln(n)} \Rightarrow K^2 \geq 8c K \ln(n) \\ K \geq 8c \ln(n)$$

$$\begin{aligned} \text{then } &\geq \frac{1}{4} \text{ nbrs red} \\ &\geq \frac{1}{4} \text{ nbrs blue} \end{aligned}$$

W.h.p: $\Theta(k \log n)$ gaps
 $\geq (\frac{1}{2} - \delta)$ red/blue nbs if $k \geq \Theta(\log n)$

Other problems

- ① Find median of array
- ② Find average of array ??

} Similar??

Random Graphs

Build $G=(V,E)$ as follows:

① $V = \text{set of nodes}$, $E = \emptyset$

② For each $u \in V$: add $2 \log(n)$ random edges
 $\rightarrow (u, ?)$

Claim: G is connected [similar to pset]

Claim: G has diameter $O(\log n)$ w.h.p.

Analysis:

\rightarrow Principle of deferred decision

\rightarrow start at $u \rightarrow$ choose edges $(u, ?)$

\rightarrow find nodes in N_1

\rightarrow choose edges for $N_1 \Rightarrow$ find nodes N_2

etc.

Notation:

Let $N_0 = \{u\} \leftarrow$ start with any u

Let $N_1 = \text{nbrs on edges chosen by } N_0$

Let $N_{j+1} = \text{nbrs on edges chosen by } N_j \searrow \bigcup_{i=1}^j N_i$

How big is N_1 ? $2C \log(n)$

Now big is N_2 ?

→ Choose $|N_1| C \log n$ edges

→ Let $X_j = 1$ if edge hits new node
= 0 otherwise

$$\rightarrow \Pr(X_j) = \frac{|N_1| + |N_0|}{n} \leq \frac{2C \log n}{n} \leq \frac{1}{2}$$

$$E[|N_2|] = E\left[\sum X_j\right] = |N_1| C \log n \cdot \frac{1}{2} = \mu \geq C \log n$$

$$C > 16$$

$$\Pr(|N_2| \leq \left(\frac{1}{2}\right)\mu) \leq e^{-\mu \left(\frac{1}{2}\right)^{\frac{1}{2}}} \leq e^{-2C \log n} \leq \frac{1}{n^2}$$

$$\Rightarrow \text{w.p.} \geq 1 - \frac{1}{n^2}, |N_2| \geq \log(n) \cdot |N_1| \geq 2|N_2|$$

For each node: divide edges selected into 2 parts:

Phase 1) $C \log n$

Phase 2) $C \log n$

⇒ Use Phase 1 edges to grow set

⇒ Use Phase 2 edges to reach the rest

Let $T_j = N_0 \cup N_1 \cup \dots \cup N_{j-1} \leftarrow$ all nodes found
in first $j-1$ steps

Phase 1) Expansion: $T_j \leq \frac{n}{2}$

Consider $N_j: T_j \leq \frac{n}{2}$

\rightarrow Choose edges for N_{j-1}

$\rightarrow \geq |N_{j-1}| \cdot \frac{c}{2} \log n$ edges

$\rightarrow X_i = 1$ if edge i hits new node (not in T_j)

$\rightarrow \Pr(X_i) \geq 1 - \frac{|T_j|}{n} \geq \frac{1}{2}$

$$M = E[|N_j|] \geq E[\sum X_i] \geq |N_{j-1}| \frac{c}{2} \log n \geq \frac{1}{6} \log n$$

$$\Pr[|N_j| \leq \frac{1}{2} M] \leq e^{-\frac{M}{2}} \leq e^{-\frac{1}{12} \log n} \leq \frac{1}{n^2}$$

$$\Rightarrow \text{w.p.} \geq 1 - \frac{1}{n^2}, \quad |N_j| \geq \frac{1}{2} M$$

[Why don't we condition on
 $|N_{j-1}|$??]

$$\Pr[\text{any } N_j \text{ fails}] \leq \frac{\log n}{n^2}$$

$$\Rightarrow \text{for } j^* = \log(n), \text{ either } T_{j^*} \geq \frac{n}{2} \text{ or } |N_{j^*}| \geq 2^{\log n} |N_0| \geq n$$

$$\text{End of phase 1) } T_{j^*} \geq \frac{n}{2}$$

Phase 2) Finish

Let E' be the set of edges chosen by nodes in T_{j^*}

$$|E'| \geq \frac{n}{2} \cdot c \log n$$

$$\text{What is } \Pr(V \text{ not hit by } E') \leq \left(1 - \frac{1}{n}\right)^{\frac{n}{2} c \log n} \leq \frac{1}{n^3}$$

$$\Pr(\text{any node in } V \text{ not hit}) \leq n \cdot \frac{1}{n^3} \leq \frac{1}{n^2}$$

$$\Pr(\text{Phase 1 fail or Phase 2 fail}) \leq \frac{\log n}{n^2} + \frac{1}{n^2} \leq \frac{1}{n}$$

Question: Bound \max degree of graph.