

15 → 19

CS473 Algorithms: Lecture 6 (2022-03-22)

logos:

last lecture: randomized algo

today: randomized algo

- closer pair - $O(n \log n)$ deterministic Δ lecture 2 divide and conquer
 - verification Δ insert points
 - randomly order points Δ create local subproblems use divide and conquer
 Δ the prob of many updates Δ

Q: what is complexity of randomized algo?

Δ we've been doing this Δ we've been using expected runtime

runtime is a random variable

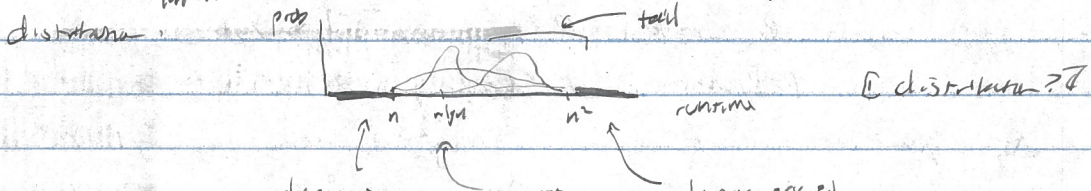
has distribution Δ how to quantify? Δ

recall quicksort

given $a = a_1, \dots, a_n$ integers, want to sort

$T(a)$ = runtime of quicksort on a

$$f(n) = \max_{a: |a|=n} \mathbb{E}[T(a)] = O(n \log n)$$



Q: $P[T(a) \geq \frac{3}{2} n^2] \leq \frac{1}{10}$ Δ really bad runtime for sorting Δ always $\geq n$ Δ expression Δ always $\leq O(n^2)$
 $\leq \frac{1}{2} n$? Δ also half of time extreme Δ no limit?
 $\leq \frac{1}{2} \sqrt{n}$? Δ small prob, acceptable Δ yes [will see?]
 $\leq \frac{1}{2} \sqrt{\log n}$? Δ so small, essentially 0 Δ yes [will see?]

def: a randomized algo runs in time $T(n)$ with probability $1 - \delta(n)$ if

for all inputs of size n , the algorithm runs in time $T(n)$, and outputs the correct answer with probability $\geq 1 - \delta(n)$

Algorithm runs in time $T(n)$ with high probability if $\delta(n) = \frac{1}{n}$ Δ error $\rightarrow 0$ quickly Δ correctness not random Δ two types of randomness

mk: quicksort will always eventually output the correct answer, the runtime is random next lecture will see algorithm whose runtime is not random, correctness only with high probability

Q: $P[X \geq \mathbb{E}[X]] \geq \frac{1}{2}$? Δ we've been using $\mathbb{E}[X]$ as measure of runtime Δ narrow range that is meaningful, as $X \geq \mathbb{E}[X]$ in our situation Δ

lem [Markov's inequality]: $X \geq 0$ and var Δ non-negative & avg: Δ

$$P[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

$$\mathbb{E}[X] = \underbrace{\mathbb{E}[X | X \geq a]}_{\geq a} \cdot P[X \geq a] + \underbrace{\mathbb{E}[X | X < a]}_{\leq a} \cdot P[X < a]$$

$$\geq a \cdot P[X \geq a]$$

cor: $X \geq 0, k \geq 0. P[X \geq k \cdot \mathbb{E}[X]] \leq \frac{1}{k}$ Δ can't deviate too far from expectation Δ runtime is non-negative! Δ quicksort

cor: $P[T(a) \geq k \cdot O(n \log n)] \leq \frac{1}{k}$ Δ goes to 0 as quickly Δ $\geq \frac{13}{10}$ Δ $O(n \log n)$

Q: do larger?

algo: quicksort - recursive - while

- run $b = \text{quicksort}(a)$ for $2 \cdot \Theta(n \lg n)$ steps
 ↳ keep track of array & array size
 ↳ if b sorted, return b
 ↳ extra $\lfloor \log n \rfloor$ steps
 ↳ expected runtime
 ↳ unsorted
 ↳ constant

prop: quicksort - recursive is always correct (check)

prop: quicksort - recursive runs in $O(k n \lg n)$ time with probability $1 - Y_2^k$ [vs $1 - Y_2^k$]

pf: check - checking if b is sorted is $O(n)$ time & simple pass, look for violations

let L be number of loops

Chn: runtime is $L \cdot (2 \cdot \Theta(k n \lg n) + O(n))$ (check)

$X_i = \begin{cases} 1 & \text{if attempt } i \text{ of quicksort succeeds in } 2 \cdot \Theta(k n \lg n) \text{ time} \\ 0 & \text{else} \end{cases}$ [has many loops?]

Chn = $P[X_i = 0] = P[\text{quicksort runtime} \geq 2 \cdot \Theta(k n \lg n)] \leq Y_2$

co = $P[L > k] = P[X_1 = \dots = X_k = 0] \leq Y_2^k$

$\Rightarrow P[\text{runtime} > k \cdot \Theta(k n \lg n)] \leq Y_2^k$

co: quicksort - recursive runs in $O(n \lg^2 n)$ time w.p. $1 - Y_2$ [why?]

↳ can generalize to generic way so turn expected runtime into w.p. also
 ↳ slight penalty in runtime
 $O(n \lg^2 n)$ time w.p. $1 - Y_2^c$, any $c = \Theta(1)$
 $\geq n^2/10$ time w.p. $Y_2^{\Theta(n \lg n)}$ [tradeoff runtime and success prob]

Q: avoid penalty turning expected runtime into w.p. runtime? [better probabilistic tools] [better analysis of algo]

fact: quicksort runs in $O(n \lg n)$ time with probability $1 - Y_2^c$, any $c = \Theta(1)$

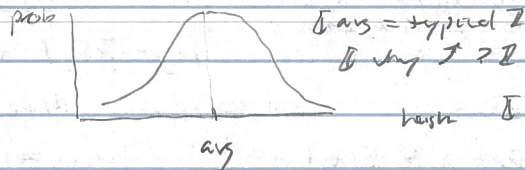
$\geq n^2/10$ time w.p. $Y_2^{\Theta(n \lg n)}$

↳ better probabilistic tools

Q: tail bounds?

Q: what is expectation? [why should we expect expectation?]

eg: distribution of human height



Thm: Central Limit Theorem

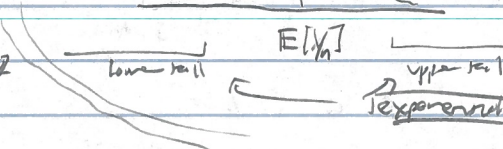
X_1, \dots, X_n random var over \mathbb{R} , identical and independent, $E[X_i] < \infty$

then $Y_n = \frac{X_1 + \dots + X_n}{\sqrt{n}} \rightarrow$ gaussian distribution as $n \rightarrow \infty$



Q: (Gaussian) kernel?

↳ needed to also analyze



upper tail exponentially small & much better than Markov

thm [Chernoff bound]: [family of bounds]

$X_1, \dots, X_n \in \{0, 1\}$ ^{binary & independent}, $E[X_i] = p_i$ ^{may not be ident}

$$X = X_1 + \dots + X_n \quad \delta \geq 0 \quad \Pr[X \geq (1+\delta) E[X]] \leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}} \right)^{E[X]} \quad \text{messy}$$

$$\epsilon \geq 0 \quad \Pr[|X - E[X]| \geq \epsilon \cdot n] \leq 2 \cdot e^{-\epsilon^2 n / 4} \quad \text{less general}$$

[do applications, then prove, rest on paper]

eg (polling)

$S \subseteq \mathcal{P}$ estimate $\frac{|S|}{|\mathcal{P}|}$ [fraction of voters for candidate]

also: pick $\sigma_1, \dots, \sigma_t \in \mathcal{P}$ w/ random [t parameters]

return $\sum_{i=1}^t \mathbb{1}[\sigma_i \in S]$

complexity: $O(t)$ [what is +?] ^{if}

concentrate - $X_i = \mathbb{1}[\sigma_i \in S]$ $E[X_i] = \frac{|S|}{|\mathcal{P}|} = \mu$

$X = \sum_{i=1}^t X_i$ $E[X] = \mu \cdot t$

$Y = \frac{1}{t} X$ $E[Y] = \mu$

$\Pr\left[\frac{1}{t} |X - E[X]| \geq \epsilon\right] \leq 2e^{-\epsilon^2 t / 4} \leq .05 \quad \text{if } t \geq 8 \cdot 738 \quad \text{if } \epsilon = .05$

$\frac{1}{t} |X - E[X]|$

if $t \geq 8 \cdot 738$ [pretty easy from before]

eg (balls into bins)

m balls, n bins, each ball thrown into random bin. take $\lfloor mn \rfloor$

Q: what is the max load of bin j ? [interesting way]

$X_{i,j} = \mathbb{1}[\text{ball } i \rightarrow \text{bin } j]$ $E[X_{i,j}] = 1/n$

$Y_j = \sum_{i=1}^m X_{i,j}$ = # balls in bin j $E[Y_j] = m \cdot 1/n = 1$ [same as hashing]

$\Pr[Y_j \geq c] \leq \left(\frac{e^{c-1}}{c}\right)^1 = e^{-\Theta(c \ln c)}$ [small!]

Q: what is the max load over $\lfloor n \rfloor$ bins?

$\Pr[\max_j Y_j \geq c] = \Pr[\exists j Y_j \geq c]$

$\leq \sum_j \Pr[Y_j \geq c]$

$= n \cdot e^{-\Theta(c \ln c)}$

choose $c = k \cdot \frac{\ln n}{\ln \ln n}$

$= n \cdot \frac{1}{n^{\Theta(k)}} = \frac{1}{n^{\Theta(k)}}$ [why max $\leq O(\frac{\ln n}{\ln \ln n})$]

$E[\max_j Y_j] = E[Y | Y < k \frac{\ln n}{\ln \ln n}] \cdot \Pr[\dots] + E[Y | Y \geq k \frac{\ln n}{\ln \ln n}] \cdot \Pr[\dots]$

$\leq k \cdot \frac{\ln n}{\ln \ln n} \leq 2 \quad \leq n \cdot \frac{1}{n^{\Theta(k)}} \leq \frac{1}{n^{\Theta(k)}}$

$= O\left(\frac{\ln n}{\ln \ln n}\right)$ for $k = \Theta(1)$

[each bin has $E=1$ individually, but max load is expected to be higher]

thm [Chernoff]. $X_1, \dots, X_n \in \Sigma_{0,1}$ random variables, independent, $\mathbb{E}[X_i] = p_i \in [0,1]$

$X = X_1 + \dots + X_n, \sigma > 0, \Pr[X \geq (1+\delta)\mathbb{E}[X]] \leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right) \mathbb{E}[X]$

idea = Markov's inequality $\left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right) \mathbb{E}[X]$
 so [transform] variable

$t > 0$ parameter, $x \mapsto e^{tx}$ is strictly increasing function

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

$\mathbb{E}[e^{tX}] = \mathbb{E}[e^{t(X_1 + \dots + X_n)}] = \mathbb{E}[e^{tX_1 + \dots + tX_n}]$
 $= \mathbb{E}[e^{tX_1} \dots e^{tX_n}]$ [independent]

$\mathbb{E}[e^{tX_i}] = \sum_{X_i=1} e^t \cdot p_i + \sum_{X_i=0} e^0 \cdot (1-p_i) = 1 + p_i(e^t - 1) \leq e^{p_i(e^t - 1)}$
 $\leq \prod_i e^{p_i(e^t - 1)} = e^{\sum p_i(e^t - 1)} = e^{\mu(e^t - 1)}$

$\leq \frac{e^{\mu(e^t - 1)}}{e^{t(1+\delta)\mu}}$
 choose $t = \ln(1+\delta)$
 $= \frac{e^{\mu((1+\delta) - 1)}}{(1+\delta)^{(1+\delta)\mu}} = \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^\mu$

- today: randomized algo
- complexity of randomized algo
 - expectation & probability
 - Chernoff bound
 - tails & more advanced
 - tail bounds
 - Markov
 - quicksort - worst case (better if tail)
 - CLT
 - Chernoff
 - padding
 - balls - n - bins

next: lemma = randomized algo

logistics = pres 6 due F17