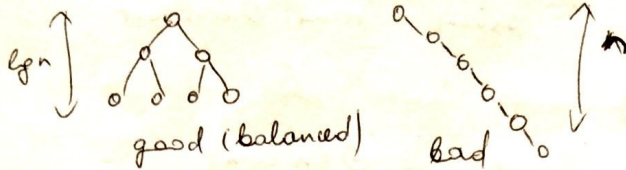


6.046
Lecture 9

Relation of BSTs to Quicksort
Analysis of Random BST

①

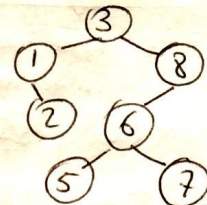


Randomly Built BSTs

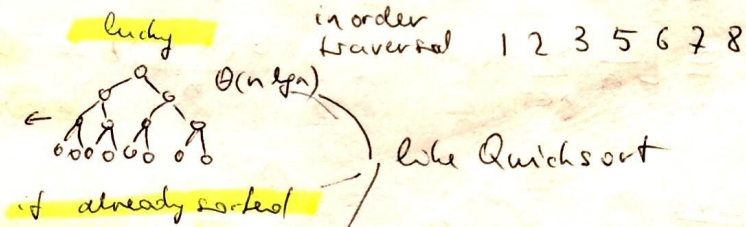
Ex1 $A = [3 | 1 | 8 | 2 | 6 | 7 | 5]$

BST sort (A)

$T \leftarrow \emptyset$
for $i \leftarrow 1$ to n
do Tree-Insert ($T, A[i]$)
Inorder-Tree-Walk (root[T])



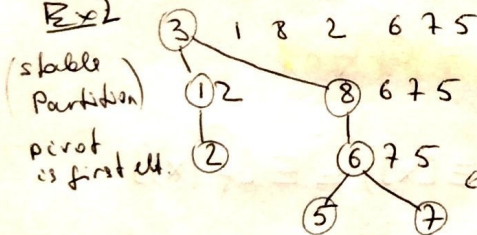
Time $O(n)$ for walk
 n Tree-Inserts
to $O(n \lg n)$
all the time
 $\rightarrow O(n^2)$



Relation to Quicksort

Comparisons that
BST sort makes are exactly
the same as comparisons that quicksort makes, though in a different order.

Ex2



same tree
as in Ex1
same comparisons as in
Ex1, but in a different order

$$\text{Time} = \sum_{x=1}^n \text{depth}(x)$$

Randomized BST Sort

- ① Randomly permute A
- ② BST sort (A)

← equivalent to picking random elt as pivot in rand. Quicksort

Time = time (rand. Quicksort)

$$E[\text{Time}] = E[\text{time (rand. Quicksort)}] = \Theta(n \lg n)$$

Randomly built BST

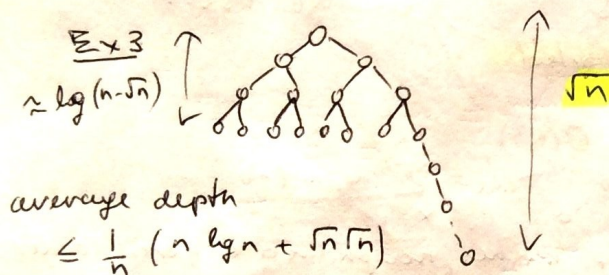
= tree resulting from randomized BST sort, without the in-order traversal

$$\text{Time (BST sort)} = \sum_{x \in T} \text{depth}(x) \quad \leftarrow \text{random variables}$$

$$E[\text{Time (BST sort)}] = \Theta(n \lg n)$$

$$E\left[\frac{1}{n} \sum_{x \in T} \text{depth}(x)\right] = \frac{\Theta(n \lg n)}{n} = \Theta(\lg n)$$

const. average depth in the tree



$$\text{av. depth} = O(\lg n), \text{ height} = \sqrt{n}$$

knowing that the average depth is $\Theta(\lg n)$

\Rightarrow height is $O(\lg n)$

Theorem $E[\text{height of rand. built BST}] = O(\lg n)$

Proof outline:

- ① Prove Jensen's inequality: $f(E[X]) \leq E[f(X)]$ for convex function f
- ② Instead of analyzing $X_n = \text{v.v. of height of BST on } n \text{ nodes}$, analyze $Y_n = 2^{X_n}$

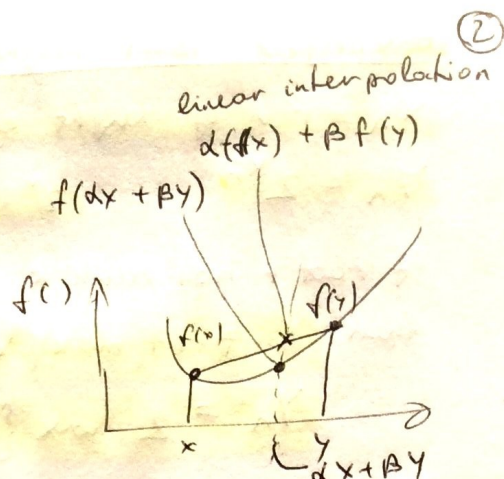
③ Prove that $E[Y_n] = O(n^3)$

④ Conclude that

$$2^{E[X_n]} \leq E[2^{X_n}] = E[Y_n] = O(n^3)$$

$$\Rightarrow E[X_n] \leq \lg O(n^3) = 3 \lg n + O(1)$$

① $f: \mathbb{R} \rightarrow \mathbb{R}$ is convex if
for all $x, y \in \mathbb{R}$
and all $\alpha, \beta \geq 0, \alpha + \beta = 1$
 $f(\alpha x + \beta y) \leq \alpha f(x) + \beta f(y)$



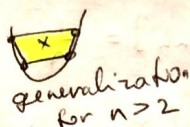
Lemma: if $f: \mathbb{R} \rightarrow \mathbb{R}$ is convex,

and $x_1, \dots, x_n \in \mathbb{R}$

and $\alpha_1, \dots, \alpha_n \geq 0$

and $\sum_{k=1}^n \alpha_k = 1$

then $f\left(\sum_{k=1}^n \alpha_k x_k\right) \leq \sum_{k=1}^n \alpha_k f(x_k)$



lin. comb $\alpha p + \beta q$, a) $\alpha + \beta = 1$
only a) \rightarrow entire line b) $\alpha, \beta \geq 0$
a) and b) only the p-q segment

Proof: Induction on n

Base: $n=1$ $\alpha_1=1 \Rightarrow f(1x_1) \leq 1 f(x_1) \checkmark$

$n=2$ by def. of convexity

Ind. Step $f\left(\sum_{k=1}^n \alpha_k x_k\right) = f\left(\alpha_n x_n + (1-\alpha_n) \sum_{k=1}^{n-1} \frac{\alpha_k}{(1-\alpha_n)} x_k\right)$

by convexity

sum to 1

sum up to 1

$$\leq \alpha_n f(x_n) + (1-\alpha_n) f\left(\sum_{k=1}^{n-1} \frac{\alpha_k}{1-\alpha_n} x_k\right)$$

by IH

$$\leq \alpha_n f(x_n) + (1-\alpha_n) \sum_{k=1}^{n-1} \frac{\alpha_k}{1-\alpha_n} f(x_k)$$

$$= \sum_{k=1}^n \alpha_k f(x_k) \checkmark$$

Jensen's inequality

$f(E[X]) \leq E[f(X)]$, if f convex, X is integer r.v.

Proof: $f(E[X]) = f\left(\sum_{x=-\infty}^{\infty} x \Pr\{X=x\}\right) \leq \sum_{x=-\infty}^{\infty} \Pr\{X=x\} f(x)$

reclustering
of sum

$$= \sum_{y \in \text{range}(f)} \underbrace{\sum_{x: f(x)=y} \Pr\{X=x\}}_{\Pr\{f(X)=y\}} = E[f(X)] \checkmark$$

Expected BST height analysis

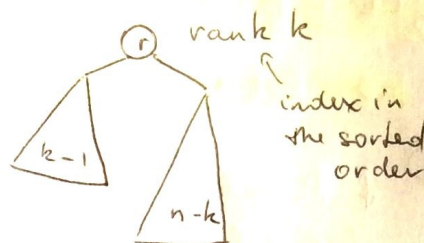
X_n = r.v. of height of randomly built BST on n nodes.

$Y_n = 2^{X_n}$ 2^x is convex

if root r has rank k

then $X_n = 1 + \max\{X_{k-1}, X_{n-k}\}$

$Y_n = 2 \max\{Y_{k-1}, Y_{n-k}\}$



better for recurrence analysis:
2. subproblem

define indicator r.v.s

$Z_{nk} = \begin{cases} 1 & \text{if root has rank } k \\ 0 & \text{otherwise} \end{cases}$

$P_r\{Z_{nk} = 1\} = E[Z_{nk}] = \frac{1}{n}$

$Y_n = \sum_{k=1}^n Z_{nk} [2 \max\{Y_{k-1}, Y_{n-k}\}]$

$E[Y_n] = E[\sum_{k=1}^n Z_{nk} [2 \max\{Y_{k-1}, Y_{n-k}\}]]$

$= \sum_{k=1}^n E[Z_{nk} [2 \max\{Y_{k-1}, Y_{n-k}\}]]$ linearity

$= 2 \sum_{k=1}^n \underbrace{E[Z_{nk}]}_{1/n} E[\max\{Y_{k-1}, Y_{n-k}\}]$ independence

$\leq \frac{2}{n} \sum_{k=1}^n E[Y_{k-1} + Y_{n-k}]$ a bit loose
linearity

$= \frac{4}{n} \sum_{k=0}^{n-1} E[Y_k]$

claim: $E[Y_n] \leq cn^3$

Proof Substitution, Base $n = \Theta(1)$, if c is sufficiently large

inductive step: $E[Y_n] \leq \frac{4}{n} \sum_{k=0}^{n-1} E[Y_k] \leq \frac{4}{n} \sum_{k=0}^{n-1} ck^3$ by IH

$\leq \frac{4c}{n} \int_0^n x^3 dx = \frac{4c}{n} \frac{n^4}{4} = cn^3 \checkmark$

approximate by integral

$\max(a, b) \leq a + b$
work
 $\max(2^a, 2^b) \leq 2^a + 2^b$
better!!

$E[X_n] \leq \lg[cn^3] = 3 \lg n + O(1)$

$E[X_n] \approx 2.9882 \cdot \lg n$ [Devroye 1986]

Very tight bound