Algorithms
Professor John Reif

## ALG 2.1

### Randomized Algorithms for Selection and Sorting:

(a) Randomized Sampling
(b) Selection by Randomized Sampling
(c) Sorting by Random
     Splitting:  Quicksort and
     Multisample Sorts

Main Reading Selections:
   CLR, Chapters 8, 10
Auxillary Reading Selections:
   AHU-Design, Sections 3.5-3.7
   BB, Sections 4.5, 4.6
   AHU-Data, Section 8.3
   Handout:  "Derivation of
   Randomized Algorithms"

# *Comparison Problems*

*input*
   set X of N distinct keys
   total ordering < over X

*Problems*

(1) for each key x $\varepsilon$ X
   $rank(x, X) = |\{x' \; \varepsilon \; X | \; x' < x\}| + 1$

(2) for each index i $\varepsilon$ {1,...,N}
   $select\,(i, X) =$ the key x $\varepsilon$ X
   where i = rank (x, X)

(3)   $sort\,(X) = (x_1, x_2, ..., x_n)$
   where $x_i =$ select (i, X)

**(1) Comparison nodes**

$$X_i < X_j$$

Yes     No

**(2)  Random Choice Nodes**

RANDOM

$prob \ \dfrac{1}{2}$     $prob \ \dfrac{1}{2}$

3

*Algorithm samplerank $_s$(x,X)*

*begin*

    Let S be a random sample of X-{x} of size s

    *output*   $1 + \dfrac{N}{s} [\, rank \,(x,S) \,-1\,]$

*end*

4

## Slide 5

*Lemma 1*

**The *expected value*
of samplerank $_s$(x,X)
is rank (x,X)**

### *proof*

Let k=rank(x,X)
For a random y $\varepsilon$ X,

$$P\ r\ o\ b\ (\ y < x\ ) \ = \ \frac{k-1}{N}$$

Hence E (rank(x,S)) = $s \cdot \dfrac{k-1}{N} + 1$

Solving for k, we get

$$r\ a\ n\ k\ (\ x\ ,X) \ = \ k \ = \ 1 + \frac{N}{s}\ E[\ r\ a\ n\ k\ (\ x\ ,S\ )-1]$$

$$= \ E(\ s\ a\ m\ p\ l\ e\ r\ a\ n\ k\ (\ x\ ,X\ )\ )$$
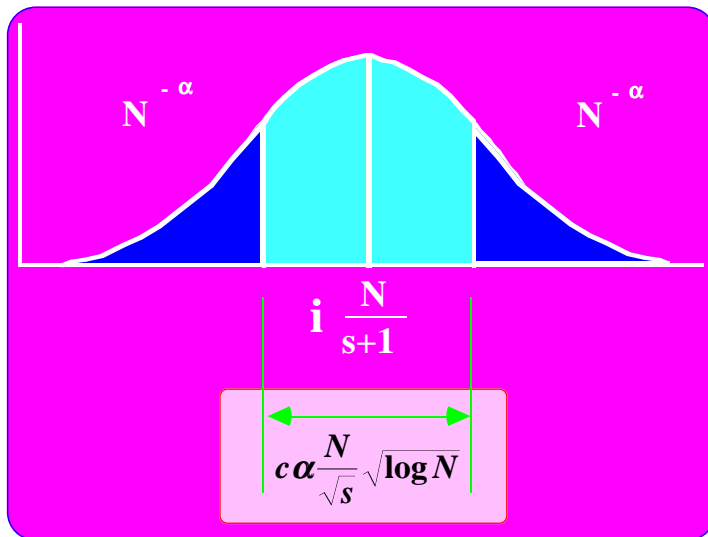
## Slide 6

**S is random sample of X of size s**

## More Precise Bounds on Randomized Sampling

**Let S be a random sampling of X**

**Let $r_i$ = rank(select (i,S),X)**

*Lemma 2*

$$\text{Prob}\left( |r_i - i\,\frac{N}{s+1}| > c\,\alpha\,\frac{N}{\sqrt{s}}\,\sqrt{\log N} \right) < N^{-\alpha}$$



$$N^{-\alpha} \qquad N^{-\alpha}$$

$$i\,\frac{N}{s+1}$$

$$c\,\alpha\,\frac{N}{\sqrt{s}}\,\sqrt{\log N}$$

*proof*

We can bound $r_i$ by a Beta distribution, implying

$$mean\ (r_i) = i\,\frac{N}{s+1}$$

$$Var\ (r_i) \leq \frac{i\,(s-i+1)}{(s+1)^2\,(s+2)}\,N^2$$

**Weak bounds follow from *Chebychev inequality***

**The Tighter bounds follow from Chernoff Bounds**

## Subdivision by Random Sampling

Let S be a random sample of X of size s

Let $k_1$, $k_2$, ... , $k_s$ be the elements of S in sorted order
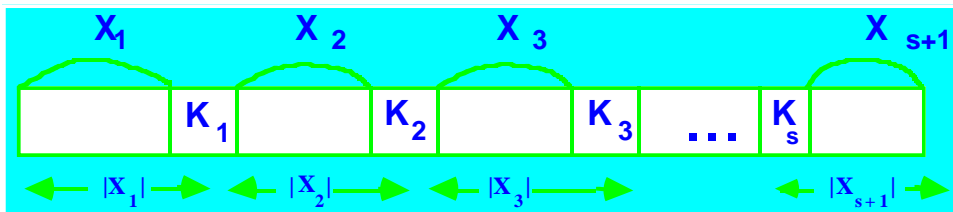
These elements *subdivide* X into

$$s+1 \text{ subsets } x_1 = \{x \ \varepsilon \ X \mid x \le k_1\}$$

$$X_2 = \{x \ \varepsilon \ X \mid k_1 < x \le k_2\}$$

$$X_3 = \{x \ \varepsilon \ X \mid k_2 < x \le k_3\}$$
$$\vdots$$
$$X_{s+1} = \{x \ \varepsilon \ X \mid x > k\}$$



## How even are these subdivisions?

### Lemma 3

If random sample S in X
is of size s and X is of size N,
then S divides X into subsets each of

$$size \le \alpha \frac{(N\text{-}1)}{s} \, ln(N) \text{ with prob} \ge 1 - N^{-\alpha}$$

## proof

**The number of** *(s+1) partitions of X* **is**

$$\binom{N\text{-}1}{s} \sim \frac{(N-1)^s}{s!}$$

**The** *number of partitions of X* **with** *one*

**block of** *size $\geq v$* **is** $\binom{N\text{-}v\text{-}1}{s} \sim \frac{(N-v-1)^s}{s!}$

**So the probability of a random (s + 1) partition having a block size $\geq$ v is**

$$\frac{\binom{N\text{-}v\text{-}1}{s}}{\binom{N\text{-}1}{s}} \sim \left(\frac{N-v-1}{N-1}\right)^s = (1-\tfrac{1}{Y})^s \ \text{for } Y = \frac{N-1}{v}$$

$$\leq (1-\tfrac{1}{Y})^{Y\left(\frac{s}{Y}\right)} \sim e^{-\frac{s}{Y}} = e^{-\frac{sv}{N-1}}$$

$$\leq N^{-\alpha} \ \text{if} \ v = \alpha\frac{(N\text{-}1)}{s} \, ln \ N$$

**since** $\left(1\text{-}\frac{1}{Y}\right)^Y < e^{-1}$

11

---

**"canonical selection algorithm"**

*Algorithm*

**can select (i,X)
input set X of N keys
index i $\varepsilon$ {1,...,N}**

**[0] if N=1 then** *output* **X**

**[1] select a bracket B of X, so that
select (i,X) $\varepsilon$ B with high prob.**

**[2] Let $i_1$ be the number of keys
less than any element of B**

**[3]** *output* **can select (i-i$_1$, B)**

*Note:* **B found by** *random sampling*
**(also must cover cases of low prob
to** *always get correct output*)

12

## Hoar's Selection Algorithm

*lgorithm*   **Hselect  (i,X)  where  1≤ i ≤N**

*begin*

    *if*   **X = {x}** *then output*   **x** *else*

    **choose a** *random splitter*   **k ε X**

    **let B = {x   ε X | x < k}**

    *if* **|B| ≥ i** *then output*   **Hselect (i,B)**

    *else output*   **Hselect (i-|B|, X-B)**

*end*

**sequential  time  bound  T(i,N)  has  mean**

$$\overline{T}\,(i,N) = N + \frac{1}{N}\left[\sum_{j=1}^{i} \overline{T}\,(i-j,\ N-j) + \sum_{j=i+1}^{N} \overline{T}\,(i,j)\right]$$

$$= 2\,N + \min\,(i,\ N\text{-}i) + o(N)$$

**random  splitter  k ε X  Hselect(i,X)   has  two  cases**

*Case |B| < i*



*Case |B| ≥ i*



*Inefficient:*  each  recursive  call  requires  *N*
*comparisons,*  but  only

**reduces  problem  size  by  average**  $\frac{1}{2}$ **N**

## *Improved Randomized Selection*
### *by Floyd and Rivest*

### *Algorithm*
### FRselect(i,X)

**begin**
    if $X = \{x\}$ then output x else
    Choose $k_1, k_2 \in X$ such that $k_1 < k_2$
    let $r_1 = rank(k_1, X)$, $r_2 = rank(k_2, X)$

    *if* $r_1 > i$ *then* FRselect(i, $\{x \in X \mid x < k_1\}$ )

*else if* $r_2 > i$ *then* FRselect($i - r_1$, $\{x \in X \mid k_1 \le x \le k_2\}$)

    *else* FRselect($i - r_2$, $\{x \in X \mid x > k_2\}$)
**end**

### *problem:*
We must choose $k_1$, $k_2$ so
that with *high likelihood,*
$$k_1 \le select(i,X) \le k_2$$

**Choose** *random sample* $S \le X$ size s

**Define:**

$$k_1 = \text{select}\left( i\frac{(s+1)}{(N+1)} - \delta, S \right)$$

$$k_2 = \text{select}\left( i\frac{(s+1)}{(N+1)} + \delta, S \right)$$

where $\delta = \left\lceil \sqrt{d\,\alpha\,s\,\log N} \right\rceil$, d=constant
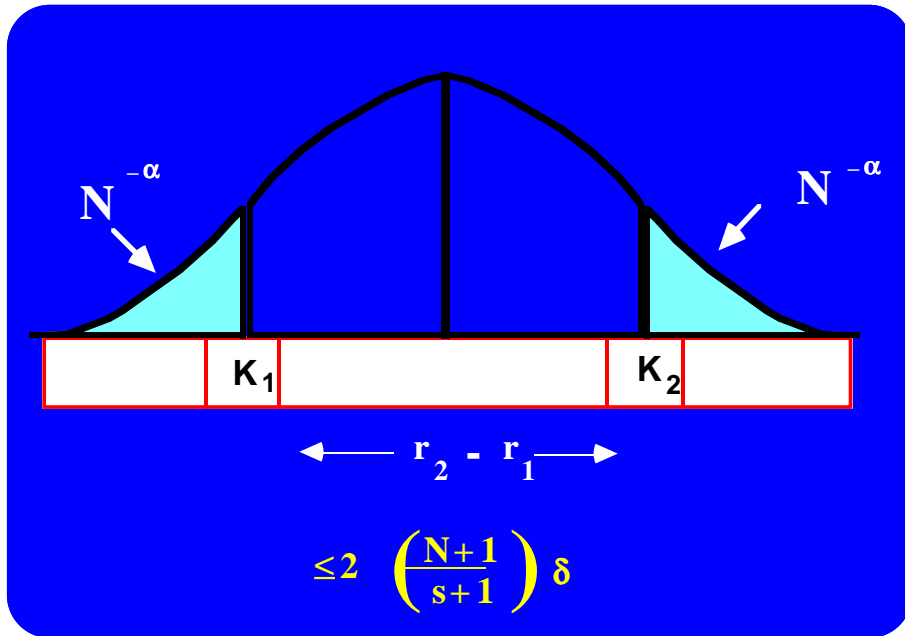
**Lemma 2 implies:**
$$Prob\ (r_1 > i) < N^{-\alpha}$$
and
$$Prob\ (r_2 < i) < N^{-\alpha}$$
where $r_1 = rank\ (k_1, X)$

$r_2 = rank\ (k_2, X)$

$$\leq 2 \left(\frac{N+1}{s+1}\right)\delta$$

## Expected Time Bound

$$\overline{T}(i,N) \leq N + 2\overline{T}(-,s)$$
$$+ \text{Prob}(r_1 > i) \cdot \overline{T}(i, r_1)$$
$$+ \text{Prob}(i > r_2) \cdot \overline{T}(i - r_1, N - r_2)$$
$$+ \text{Prob}(r_1 \leq i \leq r_2) \cdot \overline{T}(i - r_1, r_2 - r_1)$$

$$\leq N + 2\overline{T}(-,s) + 2N^{-\alpha} \cdot N + \overline{T}\left(i, 2\left\lceil \frac{N+1}{s+1}\delta \right\rceil\right)$$

$$\leq N + \min(i, N-i) + o(N)$$

$$\text{if we set } \delta < \frac{3}{\alpha} \text{ and } s = N^{\frac{2}{3}} \log N$$
$$= o(N)$$

**note**
with $\text{prob} \geq 1 - 2N^{-\alpha}$ each
*recursive call costs* only $O(s) = o(N)$
rather than N in previous algorithm

# *Randomized Sorting Algorithms*

## "canonical sorting algorithm"

**Algorithm**     **cansort(X)**
 *begin*
        *if*  x = {x}    *then output*  X  *else*
        choose a    *random sample*  S of X of size s
        *Sort S*
        *S subdivides X*    into s+1 subsets
                $X_1, X_2, ...., X_{s+1}$
    *output*   cansort($X_1$   ) cansort($X_2$   ) ... cansort($X_{s+1}$     )
 *end*

**Problem:**
    **must subdivide X into        subsets of**
            *nearly equal size*
    **to minimize number of comparisons**

        **Solution:  random sampling!**

## Hoar's Randomized Sorting Algorithm
     uses    *sample size s=1*

**Algorithm**              **quicksort(X)**
 *begin*

        *if*  |X|=1  *then output*        X else
        choose a     *random splitter*        k ε X
    *output*    quicksort({x     εX|x<k}) · (k) · quicksort({x     εX|x>k})
 *end*

### Expected Time Cost

$$\overline{T}(N) \le N - 1 + \frac{1}{N} \sum_{i=1}^{N} (\overline{T}(i-1) + \overline{T}(N-i))$$

$$\le 2 \ N \log N$$

*inefficient:*
  **need to divide problem size**
  **by $\frac{1}{2}$ with**  *high likelihood!*

## Better choice splitter is

$$k = \text{sample select}_S (\lfloor N/2 \rfloor , N)$$

**Algorithm** samplesort $_s$ (X)

begin

    *if* |X|=1 *then output* X
    choose a *random subset* S of X *size*
       s=N/log N
    k ← select( $\lfloor$s/2$\rfloor$ ,S) cost time o(N)

  *output*
     samplesort $_S$({x εX|x<k}) · (k) · samplesort $_S$({x εX|x>k})

end

---

## By Lemma 2, rank(k,X) is *very nearly the mean:*

$$\text{Prob}\left(|\text{rank}(k,X) - \frac{N}{2}| > \sqrt{d \, \alpha \, N} \, \log N\right) < N^{-\alpha}$$



**Expected Time Bounds**

$$\overline{T}(N) \leq 2\overline{T}(N_1) + N^{-\alpha} \, \overline{T}(N) \cdot N + o(N) + N - 1$$

$\approx \log (N!)$ is *optimal for comparison trees!*

## Open Problems
## in
## Selection and Sorting

**(1)** *Improve* **Randomized Algorithms**

to *exactly match lower bounds*

**on number of comparisons**

**(2)** **Can we** *de* **randomize these** **algorithms -**

**i.e., give** *deterministic algorithms*

**with the same bounds?**