

## RECAP

CHERNOFF BOUNDS for sums of Bernoulli trials (independent indicator variables)

CHERNOFF'S SOUND 1 : Let  $X_1, X_2, \dots, X_n$  be independent indicator variables (Bernoulli trials) with  $\Pr(X_i = 1) = p_i = E[X_i]$ ,  $0 < p_i < 1$ ,  $1 \leq i \leq n$ . Let  $X = \sum_{i=1}^n X_i$ ,  $\mu = E[X] = \sum_{i=1}^n p_i$

Then  $\forall \delta > 0$  :

$$\text{Prob}(X > (1+\delta)\mu) < \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^\mu$$

CHERNOFF'S SOUND 2 : Let  $X_1, X_2, \dots, X_n$  be independent indicator variables (Bernoulli trials) with  $\Pr(X_i = 1) = p_i = E[X_i]$ ,  $0 < p_i < 1$ ,  $1 \leq i \leq n$ . Let  $X = \sum_{i=1}^n X_i$ ,  $\mu = E[X] = \sum_{i=1}^n p_i$

Then  $\forall 0 < \delta < 1$  :

$$\text{Prob}(X > (1+\delta)\mu) < e^{-\delta^2\mu/3}$$

CHERNOFF'S SOUND 3 : Let  $X_1, X_2, \dots, X_n$  be independent indicator variables (Bernoulli trials) with  $\Pr(X_i = 1) = p_i = E[X_i]$ ,  $0 < p_i < 1$ ,  $1 \leq i \leq n$ . Let  $X = \sum_{i=1}^n X_i$ ,  $\mu = E[X] = \sum_{i=1}^n p_i$

Then  $\forall 0 < \delta < 1$  :

$$\text{Prob}(X < (1-\delta)\mu) < e^{-\delta^2\mu/2}$$

# ANALYSIS OF RANDOMIZED QUICKSORT

Let  $S$  be a set of  $n$  distinct integers (w.l.o.g.). We will write with

$\text{SORT}(S) = \langle x_1, x_2, \dots, x_n \rangle$  the sorted sequence of elements in  $S$ .

We can write:

ORDER STATISTICS:  $x_1 = \min(S)$   
 (O.S.)  $x_n = \max(S)$

```

QUICKSORT(S)
if |S| ≤ 1 then return <S>
y ← CHOOSE-PIVOT(S)
S1 ← {x ∈ S : x < y} }  

S2 ← {x ∈ S : x > y} }n-1
X1 ← QUICKSORT(S1)
X2 ← QUICKSORT(S2)
return <X1, y, X2>
    
```

There are various ways of implementing CHOOSE-PIVOT:

1. return  $S[1]$

(always return the first element of the unsorted sequence). If  $S$  is already sorted ( $S[1..n] = \langle x_1, x_2, \dots, x_n \rangle$ ) we obtain

$$T(n) = T(n-1) + n-1 \Rightarrow T(n) = \Theta(n^2)$$

## 2. return MEDIAN(S)

(always return  $x_{\lfloor \frac{n+1}{2} \rfloor}$ ). This is difficult to do, and requires 2 complicated algorithm runs. running in time  $C \cdot n$ , with a very high constant  $C$  [AES, 220-222]

However, we obtain:

MEDIAN PARTITION

$$T(n) \leq 2T\left(\left\lfloor \frac{M}{2} \right\rfloor\right) + (C+1) \cdot n$$

$$\Rightarrow T(n) \propto (C+1)n \log_2 n$$

High constant factor. Also, the resulting algorithm is not in place.

## 3. return RANDOM(S).

We will study this scenario in high probability.

BASIC IDEA the two scenarios 1.

and 2. for CHOOSE PIVOT are

extreme (worst-1 vs best-2)

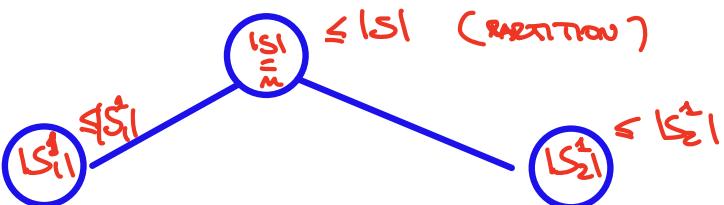
A random choice of the pivot does not guarantee a perfectly even split but will however

split  $S \setminus S^0$  into subsets  $S_1^1$  and  $S_2^1$  that are both "rather large", with "good" probability. As a consequence  $\text{YERAND}(S)$  (37) behaves more like 2 than 1.

Consider the recursion tree of quicksort!

level 0

level 1



level i:



leaves:



- (We eliminate leaves corresponding to empty cells: not all internal nodes have two children)  
 $\Rightarrow \leq n$  leaves
- For each internal node  $S_j^i$ , the sets associated to its children  $S'_j, S''_j$ :  
 $|S'_j| + |S''_j| = |S_j^i| - 1 \leq |S_j^i|$

total work at level  $i+1 \leq$  total work at level  $i$

$\Rightarrow$  At each Level  $i$ : total work at level  $i \leq n$ .

$$\Rightarrow T_{QS} \leq n \cdot (\# \text{ levels})$$

With an even split (given by  $y = \text{MEDIAN}(S)$ ):

$$\Rightarrow |S_j^i| \leq \lfloor \frac{n}{2^i} \rfloor \leq \frac{n}{2^i}$$

$$\# \text{levels} \leq \log_2 n + 1 \Rightarrow T(n) \leq (\zeta + 1) \cdot n (\log_2 n + 1)$$

### WEAKER BALANCING:

Assume that, rather than guaranteeing an even split (as given by  $y = \text{MEDIAN}(S)$ ) we are happy with a slightly more unbalanced split:  $|S_1^i|, |S_2^i| \leq \frac{3}{4}n$ .  
 $(\Rightarrow |S_1^i|, |S_2^i| \geq n - \frac{3}{4}n = \frac{1}{4}n - 1)$

We have that at level  $i$ :  $|S_j^i| \leq \left(\frac{3}{4}\right)^i n$   
At some instance  $S_j^i$ .

The maximum level  $h$  is such that

$$h = \max \{ i : n \cdot \left(\frac{3}{4}\right)^i \geq 1 \}$$

$$n \cdot \left(\frac{3}{4}\right)^i \geq 1 \Rightarrow n \geq \left(\frac{4}{3}\right)^i \Rightarrow i \leq \log_{\frac{4}{3}} n$$

$$\Rightarrow h \leq \log_{\frac{4}{3}} n$$

$$\rightarrow T(n) = O(n \cdot h) = O(n \cdot \log_{\frac{4}{3}} n) = O(n \log n)$$

The asymptotic running time stays the same even if we cannot guarantee perfect splits but only weaker balancing.

**CRUCIAL IDEA!** Weak balancing is very easy to achieve, with good probability.

**PROPERTY** Let  $SAT(S) = \langle x_1, \dots, x_n \rangle$ . If  $x_i$  is selected as a pivot, then

$$\left\lfloor \frac{n}{4} \right\rfloor + 1 \leq i \leq \lceil \frac{3}{4}n \rceil$$

then  $|S_1^+|, |S_2^+| \leq \frac{3}{4}n$

**PROOF:**  $S_1^+$  will surely contain  $x_1, \dots, x_{\lfloor \frac{n}{4} \rfloor}$   
 $S_2^+$  will contain  $x_{\lceil \frac{3}{4}n \rceil + 1}, \dots, x_n$  - Thus:

$$|S_1^+| \geq \left\lfloor \frac{n}{4} \right\rfloor, \quad |S_2^+| \geq n - \lceil \frac{3}{4}n \rceil$$

$$\begin{aligned} \text{But } |S_1^+| &= n - 1 - |S_2^+| \leq n - 1 - n + \lceil \frac{3}{4}n \rceil = \\ &= \lceil \frac{3}{4}n \rceil - 1 \leq \frac{3}{4}n \quad \lceil \frac{3}{4}n \rceil \leq \frac{3}{4}n + 1 \end{aligned}$$

$$\begin{aligned} |S_2^+| &\leq n - 1 - |S_1^+| = n - 1 - \left\lfloor \frac{n}{4} \right\rfloor \leq \\ &\stackrel{\left\lfloor \frac{n}{4} \right\rfloor = \left\lfloor \frac{n}{4} \right\rfloor - 1}{\rightarrow} \leq n - 1 - \frac{n}{4} + 1 \leq \frac{3}{4}n \end{aligned}$$

How many choices of the pivot guarantee weak balancing?

$$\lceil \frac{3}{4}n \rceil - \left\lfloor \frac{n}{4} \right\rfloor \geq \frac{3}{4}n - \frac{n}{4} = \frac{n}{2}$$

**MORALE:** If I pick the pivot at random from  $S$ , I achieve weak balancing with probability  $\geq \frac{1}{2}$

**INTUTION:** On average, given a set  $S$ , I can reduce the size of all of the subinstances to  $\leq \frac{2}{3}|S|$  in no more than two levels! Thus the maximum number of levels under 3) ( $y + \text{RANDOM}(S)$ ) is  $\approx c \log_{\frac{2}{3}} n$

We transform this intuition into a high-probability analysis using Chernoff's bound 3.

Consider the recursion tree of a call of  $\text{QUICKSORT}(S)$  with random pivot. We know that:

1. the tree has  $\leq n$  leaves  $\Rightarrow$  there are  $\leq n$  distinct root-to-leaf (r2l) paths.
2. The total work per level is  $\leq n$

Let us bound the length of a fixed r2l path in the tree:

For  $a > 1$ , let  $t = a \log_{\frac{2}{3}} n$ . Let us study

$\Pr(\text{a fixed r2l path has length} > t)$

If "the r2l path has length  $> t$ ":

$\Rightarrow$  "in the first  $t = a \log_{\frac{2}{3}} n$  nodes on the path, the algorithm has made

less than  $\log_{43} n$  pivot choices that guarantee weak balancing (or the path would end by level  $t$ !)

$$E_1 \Rightarrow E_2 \Rightarrow \Pr(E_1) < \Pr(E_2)$$

We can model this event as follows: consider  $t$  indicator variables (one per level):  $X_1, X_2, \dots, X_t$ , where  $X_i = 1$  if the pivot choice at the  $i$ -th node of the path yields weak balancing. We have argued that

$$\Pr(X_i = 1) = \frac{1}{2} \quad (\text{in fact, } \geq \frac{1}{2})$$

since the  $X_i$ 's are obtained by different calls to `RANDOM`, the  $X_i$ 's are independent (Bernoulli trials). We have to study:

$$\leq \Pr\left(X = \sum_{i=1}^t X_i < \log_{43} n\right) \quad (1)$$

Observe that  $\mu = E[X] = \sum_{i=1}^t E[X_i] = t/2$   
let's rewrite (1) in "Chernoff3" form!

Determine  $\delta$ :

$$\log_{\frac{3}{4}} n = (1-\delta)\mu = (1-\delta)\frac{t}{2} = (1-\delta)\frac{\log_{43} n}{2}$$

By fixing  $\delta = 8$ , we have  $\mu = \frac{t}{2} = 4 \log_{\frac{3}{4}} n$ ,  
therefore  $\log_{\frac{3}{4}} n = \frac{1}{4} \cdot \mu \Rightarrow (1-\delta) - \frac{1}{4} \Rightarrow \delta = \frac{3}{4}$

$$(1-\delta)\mu$$

By applying Chernoff's bound 3:

$$\Pr(X < \log_3 n) = \Pr(X < (1-\delta)\mu) < e^{-\frac{\delta^2 \mu}{2}} \quad \left| \begin{array}{l} \delta = \frac{3}{4} \\ \mu = \frac{8 \log_3 n}{3} \end{array} \right.$$

$$= e^{-\frac{(\frac{3}{4})(4 \log_3 n)}{2}} = e^{-\frac{3 \log_3 n}{3}} < e^{-\log_3 n}$$

$$= e^{-\frac{\ln n}{\ln(4/3)}} \left(e^{-\ln n}\right)^{\frac{1}{\ln(4/3)}} = \frac{1}{n^{(1/\ln(4/3))}} < \frac{1}{n^3},$$

since  $1/\ln(4/3) = 3.47\dots$

We have just proved that the probability that a fixed root-to-leaf path is longer than  $\log_3 n = 8 \log_3 n$  is  $< 1/n^3$ .

However, we have  $P \leq n$  such paths in the tree!

Let  $E_i$  be the event that the  $i$ -th path has length  $> \log_3 n$ . We have proved that  $\Pr(E_i) < \frac{1}{n^3}$

We can upper bound the probability that there is a path of length  $> \log_3 n$  as

$$\Pr\left(\bigcup_{i=1}^n E_i\right) \leq \sum_{i=1}^n \Pr(E_i) < \frac{n}{n^3} \leq \frac{1}{n^2}$$

UNION BOUND

Finally we can say that

$$\Pr(\text{Trees}(n) > 8n \log_{4/3} n) \leq$$

$$\Pr(\text{the recursion tree has } > 8 \log_{4/3} n \text{ levels}) \leq$$

$$\Pr\left(\bigcup_{i=1}^R E_i\right) \leq \frac{1}{n^2}$$

or, equivalently

$$\Pr(\text{Trees}(n) \leq 8n \log_{4/3} n) \geq 1 - \frac{1}{n^2}$$

HIGH PROBABILITY !

