

# Lower Bounds for property testing algorithms

I. Deterministic lower bounds  $\Rightarrow$  probabilistic lower bounds

a difficulty:

prop testing algs are **randomized!**  
 difficult to argue about their behavior

useful lower bnd tool:

Yao's principle:

If there is a probability distribution  $D$  on union of "positive" + "negative" elements of domain, such that any deterministic algorithm of query complexity  $\leq t$  is incorrect with prob  $\geq 1/3$  for inputs chosen according to  $D$ , then  $t$  is a lower bound on **randomized** query complexity.

So average case <sup>deterministic</sup> lower bound  $\Rightarrow$  randomized worst case lower bound  
 (principle works for all types of randomized algorithms)

why?

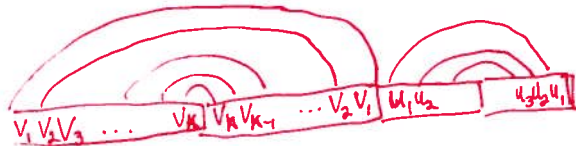
proof omitted

game theoretic view:

Alice selects deterministic alg  $A$  } payoff = cost of  $A(x)$   
 Bob selects input  $x$

Von Neuman's minimax  $\Rightarrow$  Bob has randomized strategy  
 does as well when  $A$  randomized

An example:



$$L_n = \left\{ w \mid w \text{ is } n\text{-bit string} \right. \\ \left. w = vv^R uu^R \right\}$$

concatenations of palindromes

note that testing  $L'_n = \{w \mid w = vv^R\}$  is trivial! compare  $w_n$  to  $w_1$

Thm need  $\Omega(\sqrt{n})$  queries to properly test  $L_n$   
 i.e. if  $A$  satisfies

$$\forall x \in P, \Pr[A(x) = PASS] \geq 2/3$$

$$\forall x \in \text{far from } P, \Pr[A(x) = FAIL] \geq 2/3$$

then  $A$  makes  $\Omega(\sqrt{n})$  queries

PF:

Plan: give distribution on inputs that is hard for all algorithms with  $o(\sqrt{n})$  queries.  
 Yao  $\Rightarrow$  randomized l.b. of  $\Omega(\sqrt{n})$

• wlog assume  $6/n$

• distribution on negative inputs!

should output FAIL

$N =$  random string of distance  $\geq \epsilon n$  from  $L_n$

Pf of claim 2 (idea)

To show: for every fixed set of  $o(\sqrt{n})$  queries, lots of strings in  $L_n$  follow that path.

Count # strings that agree with  $t$  queries in leaf?

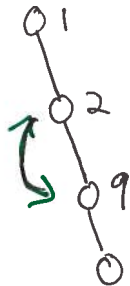
$$= 2^{n-t}$$

Count # strings in  $L_n$  that agree with  $t$  queries total?

$$\geq (2^{n-t}) - ?$$

MAIN DIFFICULTY:

must be same



Fix  $k=10$  once you see 1, that fixes what you see at 10

1	should
2	9
3	8
4	7
5	6
11	n
12	n-1



so maybe no string in  $L_n$  follows the path?



no!  $k$  could be  $\frac{n}{6} \dots \frac{n}{3}$

so for each set of queries, some  $k$ 's (but not all) are bad

• distribution on positive inputs:

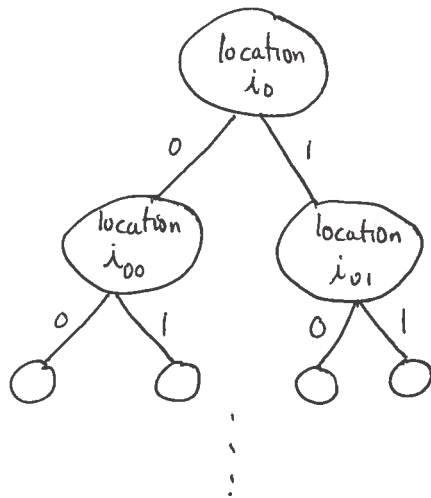
- $$P = \begin{cases} 1. & \text{pick } k \in_k \left[ \frac{n}{6}, \frac{n}{3} \right] \\ 2. & \text{pick random } v, u \text{ st. } |v|=k \\ & |u| = \frac{n-k}{2} \\ 3. & \text{output } v^R u u^R \end{cases}$$

should output Pass

an issue:  
some strings can be generated via  $\geq 1$   $k$

- distribution  $D =$   $\left\{ \begin{array}{l} \cdot \text{ flip coin} \\ \cdot \text{ if } H \text{ output according to } N \\ \text{else " " " " } P \end{array} \right.$

• Assume deterministic algorithm  $A$  has behavior above + uses  $\leq t = o(\sqrt{n})$  queries



depth  $t$ ,  $\leq 2^t$  root-leaf paths

wlog all leaves have depth  $t$

leaves labelled with A's answer following that path + seeing those bits

Note: we can calculate prob of reaching a leaf since we know input distribution



↑  
if a input reaches here, hopefully it is a "FAIL" input?

For each leaf  $l$ :

$$E^-(l) = \{ \overset{\text{inputs}}{w} \in \{0,1\}^n \mid \underbrace{\text{dist}(w, L)}_{w \text{ should fail}} \geq \epsilon n \text{ and } w \text{ reaches leaf } l \}$$

$$E^+(l) = \{ \overset{\text{inputs}}{w} \in \{0,1\}^n \cap L \mid \underbrace{w \text{ reaches leaf } l}_{w \text{ should Pass}} \}$$

each leaf  $l$  is either passing or failing, not both

Total error of  $A$  on  $D$

$$= \sum_{\substack{l \\ \text{passing}}} \Pr_{w \in D} [w \in E^-(l)] + \sum_{\substack{l \\ \text{failing}}} \Pr_{w \in D} [w \in E^+(l)]$$

↑  
should FAIL
↑  
should Pass

Claim 1 if  $t = o(n)$ ,  $\forall l$  at depth  $t$

$$\Pr_D [w \in E^-(l)] \geq \left(\frac{1}{2} - o(1)\right) 2^{-t}$$

(so negative inputs show up at all leaves & should be failed)

Claim 2 if  $t = o(\sqrt{n})$ ,  $\forall l$  at depth  $t$

$$\Pr_D [w \in E^+(l)] \geq \left(\frac{1}{2} - o(1)\right) 2^{-t}$$

(so positive inputs show up at all leaves & should be passed)

but each leaf only has one label!

Putting them together to prove full theorem

error of  $A$  on  $D$

$$\begin{aligned}
 &= \sum_{l \text{ passing}} \Pr_{w \in D} [w \in E^-(l)] + \sum_{l \text{ failing}} \Pr_{w \in D} [w \in E^+(l)] \\
 &\geq \sum_{l \text{ passing}} \left(\frac{1}{2} - o(1)\right) 2^{-t} + \sum_{l \text{ failing}} \left(\frac{1}{2} - o(1)\right) 2^{-t} \\
 &\geq \frac{1}{2} - o(1) \quad \leftarrow \text{since all leaves pass or fail}
 \end{aligned}$$

Pf of Claim 1:

idea  $N$  is close to  $U$

$+ U$  ends up uniformly distributed at each leaf  $\Rightarrow \Pr_{w \in U} [w \in E^-(l)] = 2^{-n-t}$

How much does the distribution change by using  $N$  instead of  $U$ ?

$$|L_n| \leq 2^{n/2} \cdot \frac{1}{2}$$

$\uparrow$  choice of  $u, v$        $\uparrow$  choice of  $i$

# words at distance  $\leq \epsilon$ :  $2^{n/2} \cdot \frac{1}{2} \cdot \sum_{i=0}^{\epsilon n} \binom{n}{i} \leq 2^{\frac{n}{2} + 2\epsilon \log(\frac{1}{\epsilon}) n}$

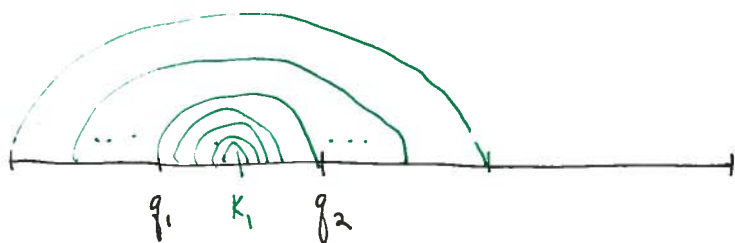
$$\text{so } E^-(l) \geq 2^{n-t} - 2^{\frac{n}{2} + 2\epsilon \log(\frac{1}{\epsilon}) n} = (1 - o(1)) 2^{n-t}$$

$\leftarrow$  # strings that follow path to leaf  
 $\leftarrow$  # words at dist  $\leq \epsilon$   
 assume  $\epsilon \ll 1/8$   
 $\epsilon$  is  $o(n)$   
 So 1st term swamps 2nd term

$$\begin{aligned}
 \text{so } \Pr_D [w \in E^-(l)] &= \frac{1}{2} \Pr_N [w \in E^-(l)] \\
 &\geq \frac{1}{2} \frac{|E^-(l)|}{2^n} \geq \left(\frac{1}{2} - o(1)\right) 2^{-t}
 \end{aligned}$$

Given leaf  $l$ , let  $Q_l \leftarrow$  indices queried along the way  
 For each of  $\binom{t}{2}$  pairs of queries  $q_1, q_2 \in Q_l$

at most 2 choices of  $k$  for which  $q_1, q_2$   
 symmetric to  $k$  or  $n/2+k$



in this case,  
 only one choice

$\Rightarrow$  # choices of  $k$  s.t.  
 no pair in  $Q_l$  symmetric  
 around  $k$  or  $n/2+k$

$$\text{is } \geq \frac{n}{6} - 2\binom{d}{2} = (1 - o(1))\frac{n}{6}$$

For these  $k$ ,  
 # strings that follow  
 path =  $2^{n/2 - |Q_l|}$

$$\begin{aligned} \text{So } \Pr_p [w \in E^+(l)] &= \sum_w \sum_k \underbrace{\Pr[w \in \mathcal{A}]}_{\frac{1}{2^{n/2}}} \cdot \underbrace{\Pr[\text{choose } k]}_{\frac{6}{n}} \cdot \mathbb{1}_{w \in E^+(l)} \\ &= \frac{1}{\frac{n}{6} \cdot 2^{n/2}} \cdot \left[ (1 - o(1)) \frac{n}{6} \right] \cdot \left[ 2^{n/2 - |Q_l|} \right] = (1 - o(1)) 2^{-|Q_l|} \\ &= (1 - o(1)) 2^{-t} \end{aligned}$$

$$\Rightarrow \Pr_p [w \in E^+(l)] = \left(\frac{1}{2} - o(1)\right) 2^{-t}$$

◻