

Lecture 12:

Testing properties of strings

Lower bounds via Yao's method

## Property Testing of Strings:

different than edit distance, which is also reasonable to consider.

def  $w$   $\epsilon$ -far from  $P_n$  if  $\forall y \in P_n$   
 $w$  &  $y$  differ on  $\geq \epsilon \cdot n$  locns

Property tester for  $P_n$ : on input  $w$

- if  $w \in P_n$ , pass (with prob  $\geq 3/4$ )
- if  $w$   $\epsilon$ -far from  $P_n$ , fail (with prob  $\geq 3/4$ )

## Palindromes

Let  $P_n = \{w \mid w \text{ is an bit string } \wedge w = vv^R\}$

Query complexity of prop tester?  $O(1)$

Algorithm:

Do  $O(1/\epsilon)$  times:

Pick random  $i$ , test if  $w_i = w_{n-i+1}$

Contrapositive:

$P \Rightarrow Q \Leftrightarrow$  equivalent

$\neg Q \Rightarrow \neg P$

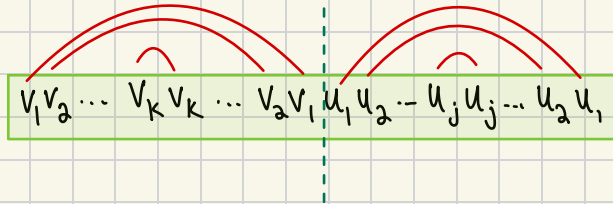
Why does it work? if  $\epsilon$ -far  $\Rightarrow$  likely to fail

equivalent to  
if not likely to fail  $\Rightarrow \epsilon$ -close

If test passes whp, then  $\leq \epsilon n$  "pairs"  $(i, n-i+1)$  don't match. Can fix each one with  $\leq \epsilon n$  changes.

## Concatenations of Palindromes

$$L_n = \{ w \mid w \text{ is } n\text{-bit string } w = v v^R u u^R \}$$



def  $w$  is  $\varepsilon$ -far from  $L_n$  if  $\forall y \in L_n$   
 $w$  &  $y$  differ on  $\geq \varepsilon \cdot n$  locns

(different than edit distance, which is also a reasonable distance to consider)

Thm if algorithm  $A$  satisfies

$$\forall x \in L_n \quad \Pr[A(x) = \text{Pass}] \geq 2/3$$

$$\forall x \text{ } \varepsilon\text{-far from } L_n \quad \Pr[A(x) = \text{Fail}] \geq 2/3$$

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

How does one prove lower bounds?

a difficulty: property testing algorithms are  
randomized

how do you argue about their behavior?

Useful tool for lower bounding randomized algorithms:

Yao's Principle:

If there is a probability distribution  $D$  on  
union of "positive" ("yes"/"pass") & "negative" ("no"/"fail")  
inputs, s.t. any deterministic algorithm of query  
complexity  $\leq t$  outputs incorrect answer  
with probability  $\geq \frac{1}{3}$  on inputs chosen according  
to  $D$ , then  $t$  is a lower bound on the  
randomized query complexity.

moral:

average case deterministic l. b.

⇔

randomized worst case l. b.

principle works for all types of randomized algorithms

why?

proof omitted

Game theoretic view:

Alice selects deterministic algorithm  $A$

Bob selects input  $x$

payoff = cost of  $A(x)$

Von Neuman's minimax  $\Rightarrow$  Bob has randomized strategy which is as good when  $A$  randomized.

# Proof of Theorem

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

without loss of generality } wlog assume  $b \mid n$

Distribution on negative inputs:

$N =$  random string of distance  $\geq \epsilon n$  from  $L_n$   
 $\leftarrow$  should output "Fail" on these

Distribution on positive inputs:

$P =$   $\left\{ \begin{array}{l} 1. \text{ pick } k \in_R \left[ \frac{n}{6} + 1, \frac{n}{3} \right] \\ 2. \text{ pick random } v, u \text{ st. } |v| = k, |u| = \frac{n-2k}{2} \\ 3. \text{ output } vv^R uu^R \end{array} \right.$

$\leftarrow$  should output "Pass" on these

note: some strings can be generated by  $\geq 1$   $k$   
e.g. 00000000...0

note:  
choice of  $k$   
 $\Rightarrow |vv^R| + |uu^R|$   
are both  $\geq n/3$   
 $\forall \in 2n/3$



Note: we can calculate probability of reaching leaf  
since we know input distribution  $\mathcal{D}$

Error of leaf  $l$ : if  $l$  is labelled:

should fail

Pass:  $E^-(l) = \{ \text{inputs } w \in \{0,1\}^n \mid w \text{ } \epsilon\text{-far + } w \text{ reaches leaf } l \}$

Fail:  $E^+(l) = \{ \text{inputs } w \in \{0,1\}^n \mid w \in L + w \text{ reaches leaf } l \}$

should pass

Total error of  $A$  on  $\mathcal{D}$ :

$$= \sum_{\substack{l \\ \text{"pass"}}} \Pr_{w \in \mathcal{D}} [w \in E^-(l)] + \sum_{\substack{l \\ \text{"fail"}}} \Pr_{w \in \mathcal{D}} [w \in E^+(l)]$$

should fail but  
reach passing leaf

should pass but reach  
failing leaf

Why is this big?

will show lots of input from both  $N$  &  $P$  end  
up at all leaves.



## Main Claims

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

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

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

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

So total error of  $A$  on  $\mathcal{D}$  is

$$= \sum_{\substack{l \\ \text{PASS}}} \left(\frac{1}{2} - o(1)\right) 2^{-t} + \sum_{\substack{l \\ \text{FAIL}}} \left(\frac{1}{2} - o(1)\right) 2^{-t} \geq \frac{1}{2} - o(1) \gg \frac{1}{3}$$

still need to prove the claims...

# Pf of Claim 1

idea:  $N$  is close to  $U$

$\& U$  would end up uniformly distributed at each leaf

$$\Rightarrow \Pr_{w \in U} [w \in E^{-}(l)] = \frac{2^{n-t}}{2^n} = 2^{-t}$$

but how much can distribution change by using  $N$  instead of  $U$ ?

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

↑ choice of  $u, v$ 
↑ choice of  $i$

# words at dist  $\leq \epsilon$  from  $L_n$ :

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

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

↑  
# strings that reach  $l$

↑  
# strings that allowed to pass ← assume worst case! they all come to leaf  $l$

$$\text{so } \Pr_w [w \in E^{-}(l)] \geq \frac{1}{2} \Pr_N [w \in E^{-}(l)] \geq \frac{1}{2} \frac{|E^{-}(l)|}{\# \text{ strings in } N} \geq \frac{1}{2} \frac{|E^{-}(l)|}{2^n} \geq \left(\frac{1}{2} - o(1)\right) 2^{-t}$$

↑ needs to come from  $N$ 
↑ upper bound on # strings in  $N$

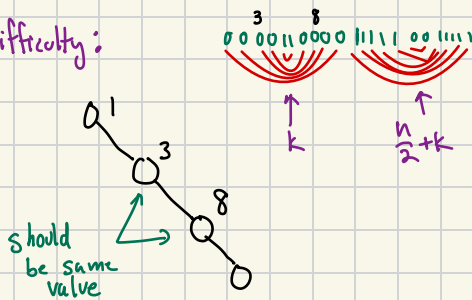
# Proof of claim 2

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

count # strings agreeing with  $t$  queries of leaf:  $2^{n-t}$

count # strings in  $L_n$  " " " " " " :  $\geq 2^{n-t} - ?$

Main difficulty:



Fix  $k=10$

Should see same value at

- |       |         |
|-------|---------|
| 1, 10 | 11, 1   |
| 2, 9  | 12, n-1 |
| 3, 8  | 13, n-2 |
| 4, 7  | ⋮       |
| 5, 6  |         |

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

😊 that's why we pick  $k$  randomly in  $[\frac{n}{6} .. \frac{n}{3}]!$   
 not all queries bad e.g. for most strings,  $3+8$  are not correlated

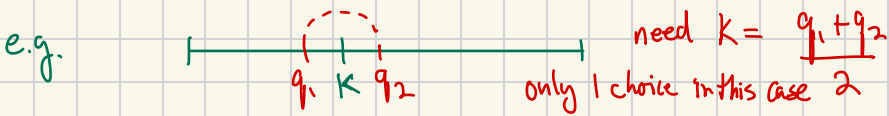
does choice of  $k$  correlate the queries?

For given leaf  $l$ , let  $Q_l \leftarrow$  indices queried on path to  $l$

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 around either  $k$  or  $\frac{n}{2} + k$



$\Rightarrow$  # choices of  $K$  s.t. no pair in  $Q_\ell$

Symmetric around  $k$  or  $\frac{n}{2} + k$  is

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

using that  $t$  is  $o(\sqrt{n})$

for these good  $K$ , # strings that follow path is  $2^{n/2 - t}$

So

$$\Pr_p [w \in E^+(Q)] = \sum_w \sum_K \Pr_p(w|K) \Pr[\text{choose } K] \cdot \mathbb{1}_{w \in E^+(Q)}$$

$$\geq \frac{1}{\binom{n}{6} (2^{n/2})} \left[ (1 - o(1)) \cdot \frac{n}{6} \right] \cdot 2^{n/2 - t}$$

$\left[ (1 - o(1)) \cdot \frac{n}{6} \right]$  ← # good  $K$

$2^{n/2 - t}$  ← # strings following path on good  $K$

$$= (1 - o(1)) \cdot 2^{-t}$$

▣