

Streaming Algs: how to handle massive data in small space

Property Testing: handle massive data in sublinear time.

Example Testing Linearity

$f: \{0,1\}^n \rightarrow \{0,1\}$ is linear if $\boxed{f(x) = \sum_{i \in S} x_i \pmod{2}}$ for some $S \subseteq [n]$

f is ϵ -far from linear if need to change $\epsilon 2^n$ entries in truth table to get linear function

Input: truth table to f

~~usually~~ usually people treat ϵ as ~~constant~~ a constant

Goal: decide if f is linear or far from linear

BLR Linearity Test

[Blum, Rubinfeld '83]

Pick $x, y \in \{0,1\}^n$ at random

Query $f(x), f(y), f(x \oplus y)$

Reject if $f(x) + f(y) \neq f(x \oplus y)$

Note: if f linear

$$f(x) + f(y) = \sum_{i \in S} x_i + \sum_{i \in S} y_i$$

$$= \sum_{i \in S} x_i \oplus y_i$$

$$= f(x \oplus y)$$

Fact: if f linear BLR always accepts

if f ϵ -far from linear BLR rejects w/prob $\geq \epsilon$.

← proof beyond scope of class. requires Fourier Analysis.

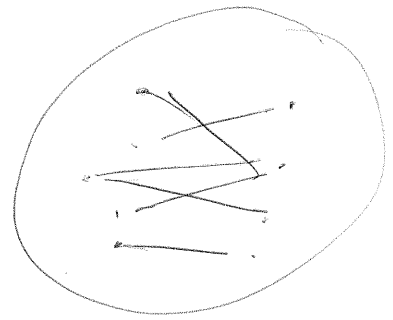
Solution: repeat BLR $O(1/\epsilon)$ times.

Accept if all tests pass.

Example testing bipartiteness

Input: undirected graph $G=(V,E)$

Goal: Is G bipartite or ϵ -far from bipartite?



need to delete ϵn^2 edges to get bipartite graph

Bipartite

Property Tester:

- ① Pick random $S \subseteq V$ $|S|=k$
- ② Query subgraph G_S : all $(e,v): v \in S$
- ③ Accept if G_S is bipartite
Reject otherwise.

Theorem

- Bipartite-Test accepts all bipartite graphs
- if $k = \tilde{O}(1/\epsilon^4)$ Bipartite-Test rejects ~~graphs~~ ϵ -far graphs w/prob $\frac{2}{3}$.

Property Testing Formal Model

Given domain D (D = functions, graphs, ...)

A property \mathbb{P} is a subset $P \subseteq D$.

Defn Given metric $\text{dist}(a,b)$ on D , $\text{dist}(a,P) := \min_{b \in P} \text{dist}(a,b)$
 a is ϵ -far from P if $\text{dist}(a,P) \geq \epsilon$

Defn T is a property tester for P if

- (1) if $a \in P$ then T accepts w/prob $\frac{2}{3}$ (completeness)
- (2) if a ϵ -far from P then T rejects w/prob $\frac{2}{3}$ (soundness)

$Q(P), Q^{NA}(P)$:
min #queries needed
for testers w/one-sided error,
nonadaptive testers

Query Complexity $Q(P) = \min \# \text{queries needed to test } P$

- (3) T has one-sided error if it accepts all $a \in P$ (i.e. w/prob 1)
- (4) T is nonadaptive if all queries chosen in advance

Lower Bounds in Property Testing

How do you prove lower bounds on query complexity?

The old way: Yao's Lemma + PAIR

The new way: Reduce from Communication Complexity

[Blais, Brody, Matulef '11]

example k-linearity

f is k-linear if $f(x) = \sum_{i \in S} x_i \pmod{2}$ for some $S \subseteq [n]$ $|S| = k$

old bounds: $O(k \log k)$

$\Omega(k)$, $\Omega(k)$ for nonadaptive testers

new bound: $\Omega(k)$

proof reduce from k-DIST:

Alice and Bob have sets $A, B \subseteq [n]$

~~with $|A| = |B| = k$~~ $|A| = |B| = k$

want to check if $A \cap B \stackrel{?}{=} \emptyset$

Fact: $R(k\text{-DIST}) = \Theta(k)$ [Hastad-Wigderson '07]

Let T be a tester for k-linearity.

k-DIST protocol:

① Alice takes A, creates $f_A: \{0,1\}^n \rightarrow \{0,1\}$

② Bob takes B, creates $g_B(x) := \sum_{i \in B} x_i$

③ Players use T to test

$$h := f_A \oplus g_B$$

For (2k)-linearity

Emulating T:

Generate query x using public randomness

Alice \rightarrow Bob $f(x)$

Bob sends $g(x)$

\Rightarrow Both players compute $h(x)$

use public randomness to

generate next query

$$f_A(x) = \sum_{i \in A} x_i$$

Again, all addition is in \mathbb{F}_2
i.e., modulo 2

$$A = \{1, 5, 7\} \quad f_A(x) = x_1 \oplus x_5 \oplus x_7$$

$$B = \{3, 5, 6\} \quad g_B(x) = x_3 \oplus x_5 \oplus x_6$$

$$h(x) = (x_1 + x_5 + x_7) + (x_3 + x_5 + x_6)$$

$$= x_1 + x_3 + x_5 + x_6 + x_7$$

Fact (1) if A, B disjoint then h is $(2k)$ -linear

(2) if A, B intersect then h is $\frac{1}{2}(k')$ -linear for some $k' < 2k$

Claim Let $k' \neq k$. Then k' -linear functions are $\frac{1}{2}$ -for from k -linear functions

proof Let $S, S' \subseteq [n]$ $|S|=k, |S'|=k'$ be given

$S \neq S'$ so $\exists i \in S \setminus S'$ (or vice versa)

For any $x \in \{0, 1\}^n$, let $x^{(i)}$ be x w/ i th bit flipped

$$\begin{aligned} x &= 011001 \\ x^{(i)} &= 011101 \end{aligned}$$

$$i \in S \text{ so } f_S(x) \neq f_S(x^{(i)})$$

$$i \notin S' \text{ so } f_{S'}(x) = f_{S'}(x^{(i)})$$

\therefore either $f_S(x) \neq f_{S'}(x)$ or $f_S(x^{(i)}) \neq f_{S'}(x^{(i)})$

$\Rightarrow f_S, f_{S'}$ differ on exactly half inputs. //

What have we proved?

(1) we can use property testing alg for k -linearity to create a communication protocol for k -DISJ

(2) the cost of the communication protocol is $2 * (\# \text{queries in tester})$

(3) But $R(k\text{-DISJ}) = \Omega(k) \Rightarrow Q(k\text{-LIN}) \geq \frac{1}{2} R(k\text{-DISJ}) = \Omega(k) //$

Testing Bipartiteness possible w/ $\tilde{O}(1/\epsilon^2)$ queries

Algorithm

- 1 Pick random $S \subseteq V$ $|S| = \tilde{O}(1/\epsilon^2)$
- 2 query (u,v) for all $u,v \in S$
- 3 ACCEPT if there is valid partition
REJECT if (S_1, S_2) not bipartite \forall partitions (S_1, S_2)

Proof We'll actually analyze this alg:

- 1 Pick $\tilde{O}(1/\epsilon)$ vertices $U \subseteq V$
- 2 Pick $S = \tilde{O}(1/\epsilon^2)$ edges query. How to partition $S = S_1, S_2$?

Defn v is influential if it has $\frac{\epsilon N}{4}$ neighbors

Claim w/prob $\geq 5/6$, at least all but at most $\frac{\epsilon N}{4}$ influential vertices adjacent to U

proof Let v be influential

$$\begin{aligned} \Pr[v \text{ not adj to } U] &\leq (1 - \frac{\epsilon}{4})^{|U|} \\ &\leq e^{-\frac{|U|\epsilon}{4}} \\ &= e^{-4 \ln 2} \\ &= \frac{\epsilon}{24} \end{aligned}$$

$\text{set } |U| = \frac{4 \ln 2}{\epsilon}$

$$\Pr[v \text{ not adj to } U] \leq \frac{\epsilon}{24}$$

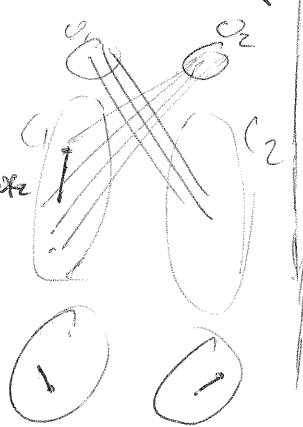
$$\Pr[\geq \frac{\epsilon N}{4} v \text{ not adj to } U] \leq \frac{1}{6} \quad (\text{Markov})$$

Now, Fix partition U_1, U_2 of U .

C_1 : neighbors of U_2

C_2 : neighbors of U_1

R_1, R_2 : a-b. partition of remaining vertices



Claim: $\leq \frac{\epsilon N^2}{2}$ violating edges touch $R = R_1 \cup R_2$

proof $\leq \frac{\epsilon N}{4}$ influential $\times N$

+ $\frac{\epsilon N}{4}$ non-influential $\times \frac{\epsilon N}{4}$

$$= \frac{\epsilon N^2}{2} \text{ total edges}$$

Claim: For any fixed partition U_1, U_2 , $\Pr[\exists$ partition S_1, S_2 that fails to witness non-bipartiteness] $\leq \frac{2}{6}$

proof There are $\geq \frac{\epsilon N^2}{2}$ violating edges in C . $|S| = \frac{16|U|}{\epsilon}$

We query ≥ 1 w/prob $\geq 1 - (1 - \frac{\epsilon}{2})^{|S|/2} \geq 1 - \frac{2^{-|S|}}{6}$

if we hit one violating pair (u,v) how to partition it? w/ say $v \in U_2$

can't put u in S_1 or S_2

if u is in S_1 , then it violates

if we put $v \in S_1$ then (u,v) witnesses when $u \in U_1$ is v 's neighbor in U_1 .

Conclusion Prob that we find any partition that doesn't witness non-bipartiteness $\leq 2^{-|U|} \cdot \frac{2^{-|U|}}{6}$

overall success: $\frac{2}{3}$

$$= \frac{1}{6}$$