

Welcome to CS590 Sublinear Algorithms

Admin.

- Brightspace ; CampusWine ; class webs. on my page.

- Grading: 4-5 HW 40%

Research project 45%

Scribe notes + peer grading HW
10%

Class participation 5% ✓

- Need scribe for today ←

- Research : • max 2 people

- See timeline : • by Feb 2 should have a topic

Schedule ←
meeting with me
before then

- Feb 16 + 2 paras
proposal

- March 23 5-10
min present

- April 29 outgoing

- April 30
written report.

Motivation for sublin. algs :

BIG DATA :

internet of things ↙
Sales transactions ↙
web pages
health data
genomic data
space discovery data
etc etc

Need algorithm design in $o(N)$

• If data can be stored but no time to read it → sublinear time algs

• If data is too big to fit in memory → sub-linear in space
if can throw some away

→ sublinear in communication
if can store it on multiple machines that communicate

Sublin-time algs

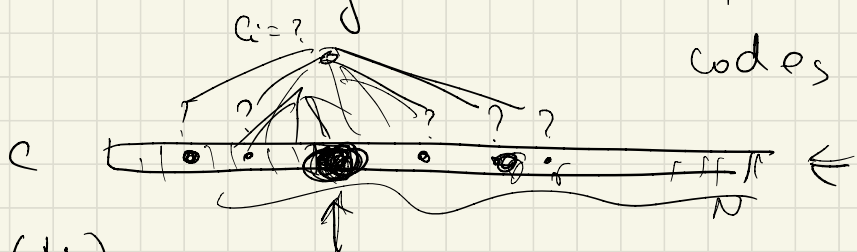
- approx algs : eg: diameter \checkmark
of connected compts -
avg degree of a graph
- property testing :
does an obj have a property or
is far from having the
property ?
(eg. is G connected or far -
is G 3-colorable or far)

Model stems from Program checking
in PL.

80's { Blum Kannan
Blum Luby Rubinfeld
formalized by Rubinfeld Sudan 90
Leads to PCP thm

Also many other local models

• Locally decodable / testable codes



$\text{poly}(N)$
 $\text{poly}(\log(N))$
 \uparrow

$c_i = ?$ if can only query few posns into word?

Type of research questions:

Can membership in specific code be tested with cl. many queries?

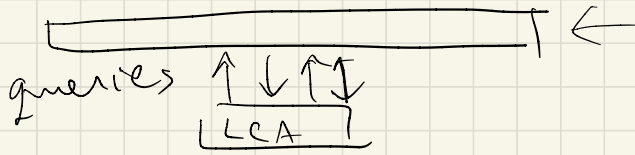
Can each entry be corrected w/hp?

Best trade offs bet : rate, distance, locality

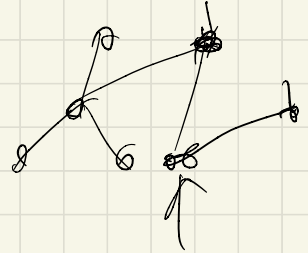
Leads to a general notion of Local Computation Algs.

Local Comput Algs (LCAs)

Rubinfeld Tamir Vardy xie '11



probe $\uparrow y_i$ $\downarrow y_i$



eg: Maximal IS.

probe node i to see if it is in the MIS selected or not.

• Based on distrib. algs where each node makes local decision & goal is to have a consistent MIS

Sublinear-space algs / streaming

Given a seq of elts appearing one at a time, with limited memory
can get eg: statistics abt the stream

• max, min, avg, median

a rand sample? ↙

estimates of # distinct elts

heavy hitters?

OR can get an approx size of

max matching

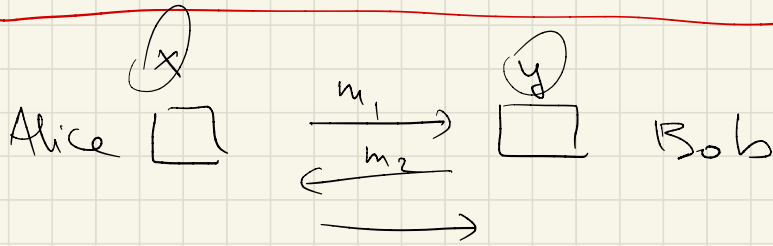
vertex cover

connected compts

avg degree?

Defined by Alon Motias Szegedy '96

Sublinear in communication



$$f(x, y)$$

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

$$f: \mathbb{F}_2 \times \mathbb{F}_2 \rightarrow \mathbb{F}_2$$

2 bits

- deterministic

- randomized

Many natural problems are hard

i.e. require almost full disclosure of inputs

Eg:

EQUALITY

$$EQ : \underbrace{\{0,1\}^N} \times \underbrace{\{0,1\}^N} \rightarrow \{0,1\}$$

$$EQ(x, y) = \begin{cases} 1 & \text{if } x=y \\ 0 & \text{otherwise} \end{cases}$$

$$\underline{\text{Det}(EQ)} \geq \underbrace{N}_{\text{Bob}}$$

INDEX

$$\text{IDX} : \underbrace{\{0,1\}^N} \times [N] \rightarrow \{0,1\}$$

$$\text{IDX}(\bar{x}, y) = \begin{matrix} x \\ y \end{matrix}$$

$$\text{Det}^{\text{one-way}}(\text{IDX}) \geq \underline{N}$$

$$\text{Rand}^{\text{one-way}}(\text{IDX}) \geq \underline{\Omega(N)}$$

SET DISJ

$$\text{DISJ} : \underbrace{\{0,1\}^N} \times \underbrace{\{0,1\}^N} \rightarrow \{0,1\}$$

$$\text{DISJ}(\bar{x}, \bar{y}) = \begin{cases} 1 & \text{if } \underline{x \cap y = \emptyset} \\ 0 & \text{otherwise} \end{cases}$$

x, y are char vectors of sets.

$$\underline{\text{Rand}(\text{DISJ})} = \underline{\underline{\Omega(N)}}$$

Obs Data stream algs \Rightarrow communic
protocol.

So lb for communication \Rightarrow
lbs for streaming. \checkmark

\downarrow

We'll probably see that

lbs for communication \Rightarrow

lbs for property testing \uparrow

• Other exs: connectivity \checkmark

• bipartiteness

• approx weight of
spanning tree \downarrow

• Recent: distributed learning of distributions \downarrow

Some basic problems

- diam. of a set of pts in \mathbb{R}^n ←
- testing if fnc is constant
- uniform sampling of a stream
- deciding connectivity of a graph in a stream.

Deterministic 2-approx of diameter

Def: D dist. metric \mathbb{R}^n

$$1) D(x, y) \geq 0 ; D(x, y) = 0 \text{ if } x = y$$

$$2) D(x, y) = D(y, x)$$

$$3) D(x, y) \leq D(x, z) + D(y, z)$$

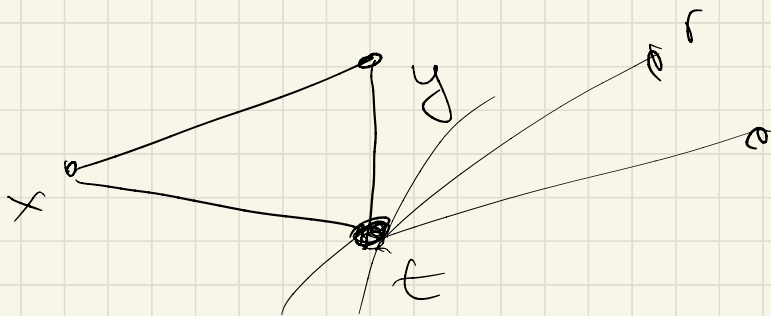
$$\text{Def: } \text{diam}(S) = \max_{x, y \in S} D(x, y)$$

Thm (Indyk): Given S, D :
∃ det. alg that
outputs d st. $\frac{\text{diam}(S)}{2} \leq d \leq \text{diam}(S)$

in time $\underline{O(|D|)}$

note: |input| = $\underline{O(|D|)}$ = $\underline{O(|S|^2)}$

|S|



$$\forall x, y \quad \forall t$$

$$D(x, y) \leq D(x, t) + D(y, t)$$

$$\leq 2 \cdot \max\{D(x, t), D(y, t)\}$$

$\frac{3}{2}$ - approx

$O(m n^{1/2})$

$$\leq 2 \max_{r \in S} \{D(r, t)\}$$

Alg: Take any $t \in S$

$$\text{output } \underline{d} = \max_{r \in S} \{D(r, t)\}$$

$$\text{Diam} \leq 2d$$

$$\leq \text{diam}(S)$$

$$\frac{\text{Diam}}{2} \leq d \leq \text{Diam} \quad \checkmark$$

Property testing

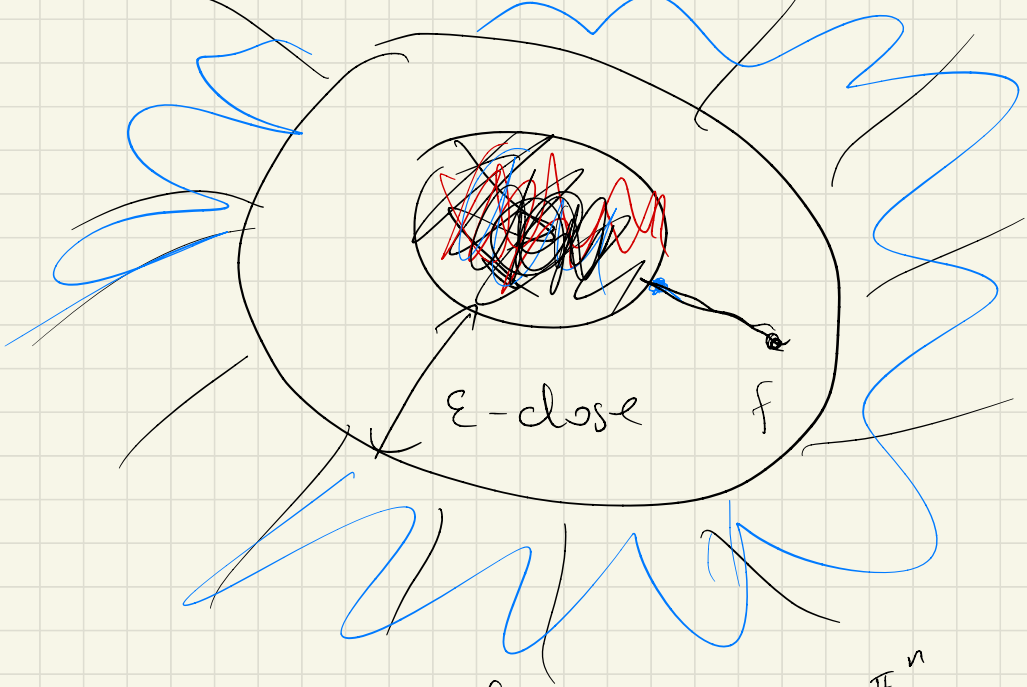
Property = a collection of obj that
all have a particular
property

$\mathcal{I}_n = \{ \text{all graph on } n\text{-vts} \\ \text{that are 3-colorable} \}$

$\mathcal{P}_n = \{ \text{all graphs on } n\text{ vts} \\ \text{that are connected} \}$

$\{ \text{---} \\ \text{that are } \Delta\text{-free} \}$

$\mathcal{P}_n = \{ \text{distrib. that are} \\ \text{bimodal} \}$



\mathbb{F}_2^n
 ϵ -far

Hamming
Rel Distance

$$f, g: D \rightarrow \mathbb{R}^k$$

$$\delta(f, g) = \frac{|\{x : f(x) \neq g(x)\}|}{|D|}$$

Dist. from f to \mathcal{P} is

$$\text{dist}(f, \mathcal{P}) = \min_{g \in \mathcal{P}} \delta(f, g)$$

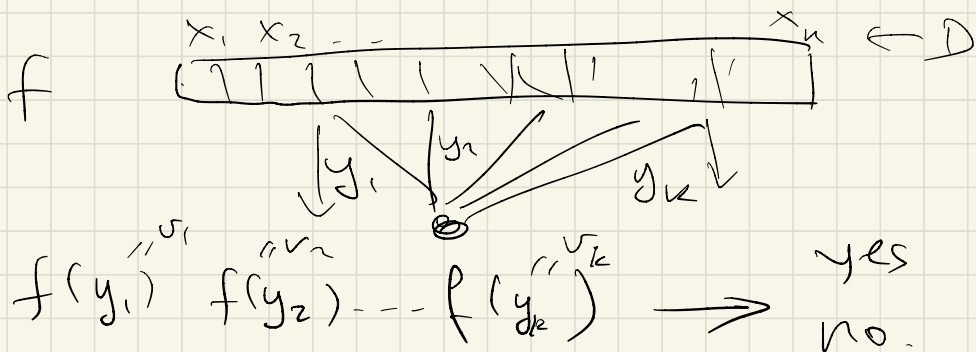
Def: \mathcal{P} is k -locally testable if \exists rand alg A with black-box access to input st .

① A makes k queries to the input.

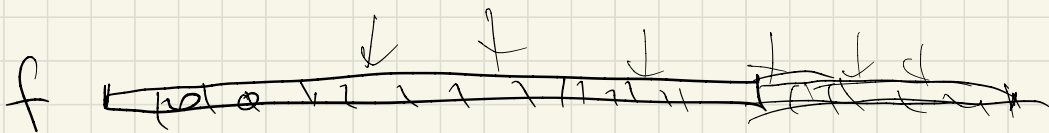
② if f is in \mathcal{P} then \leftarrow (one-sided) A accepts \leftarrow (completeness)

• if f is ϵ -far from \mathcal{P}

then $\Pr[A \text{ accepts}] < \frac{1}{3}$ (soundness)



$$P_n = \left\{ f \mid \begin{array}{l} f(x) = 1, x \in [n] \\ f: [n] \rightarrow \{0, 1\} \end{array} \right\}$$



is $f \in P_n$ or

f is ϵ -far from P_n .

(i.e. f has at least

an ϵ -frac of 0's)

$\Rightarrow \epsilon n$ values of f are 0's.

claim: P_n is $\frac{2}{\epsilon}$ -locally testable.

Disting bet $f \in P_n$

f is ϵ -far from P_n

using only $\frac{2}{\epsilon}$ many queries.
for $n \rightarrow \infty$

• Test: pick $\left(\frac{2}{\epsilon}\right)$ random x 's

$f(x_1) \quad f(x_2) \quad \dots \quad f(x_{\frac{2}{\epsilon}})$
" " " " " "
" " " " " "

if ever see 0 rej.
ow acc.

Analysis

queries: $\frac{2}{\epsilon} = O\left(\frac{1}{\epsilon}\right)$ wrt. n

completeness: $f \in P_n$ then $\Pr[\text{acc}] = 1$

soundness: f is ϵ -far from P_n

then

$$\Pr[\text{acc}] = \Pr[\text{no } 0 \text{ is hit in } \frac{2}{\epsilon} \text{ trials}]$$
$$\leq \underbrace{(1-\epsilon)^{\frac{2}{\epsilon}}} \leq \underbrace{\left(\frac{1}{e}\right)^{\frac{2}{\epsilon}}} = \frac{1}{e^2} \leq \frac{1}{2}$$