
Data Stream Algorithms in Computational Geometry

Timothy M. Chan
School of CS
Univ of Waterloo

INTRODUCTION

Data Stream: The Basic Model

makes one pass over input

uses limited (sublinear) space

Motivation

MASSIVE data sets

Brief History [Indyk]

Ancient times

(finite automata, sorting w. few passes,
median [Munro, Paterson '80])

Middle ages

Renaissance

([Alon, Matias, Szegedy '96],
[Henzinger, Raghavan, Rajagoplan '98],
& TONS of papers ...
see Muthukrishnan's survey)

Data Stream Meets CG

diameter [Feigenbaum, Kannan, Zhang '02,
Hershberger, Suri '03]

other measure problems like width

[Agarwal, Har-Peled, Varadarajan '04,
Chan (SoCG'04)]

statistical problems like range counting

[Bagchi, Chaudhary, Eppstein, Goodrich (SoCG'04)
Suri, Toth, Zhou (SoCG'04)]

clustering problems like k-median/k-means

[Har-Peled, Mazumdar (STOC'04)]

Euclidean MST/matching [Indyk (STOC'04)]

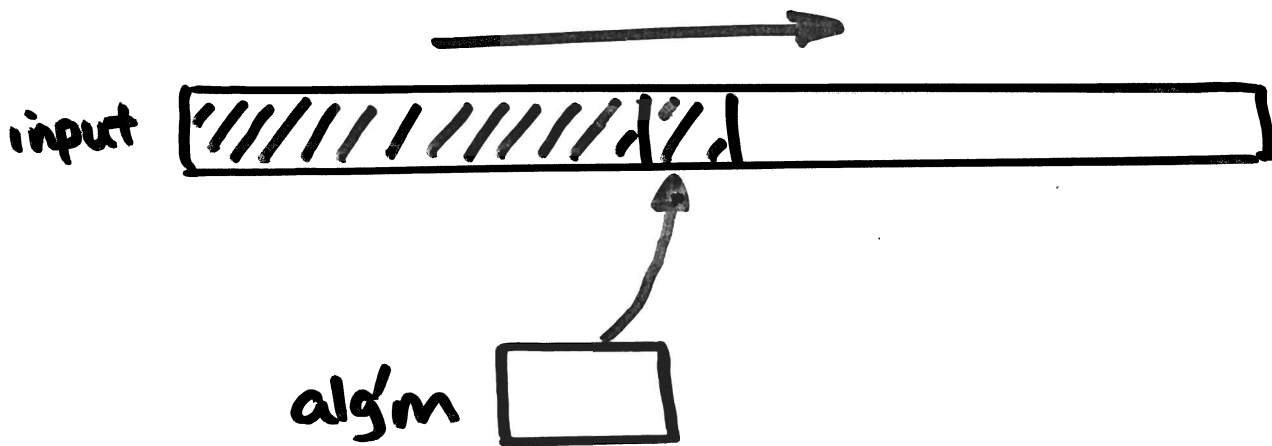
This Talk

looks at some specific examples
& some recent results

[Chan (SoCG'04),
Chan, Sadjad (ISAAC'04),
Chan, Chen (SoCG'05)]

illustrates a few techniques
& different types of
geometric streaming alg/ms

ONE-PASS ALG'S



(data structures w. sublinear space
supporting inserts)

Example 0: Approx Median in 1D

Given n pts in \mathbb{R}^1 ,

find a pt of rank $\sim \frac{n}{2} \pm \epsilon n$



Munro, Paterson's Alg'm ['80]

Standard idea: the "logarithmic" method
[Bentley, Saxe '80]

Divide set into $O(\log n)$ groups
whose sizes are distinct powers of 2



To insert:

create new group of size 1

whenever 2 groups have same size,
merge

Modified idea: "merge-and-reduce"

Replace each group with sketch

Def: RCS is a δ -sketch of S iff
ith pt of R has rank $\sim \frac{in}{r} \pm \delta n$
($n = |S|$, $r = |R|$)

Fact: \exists δ -sketch of size $1/\delta$

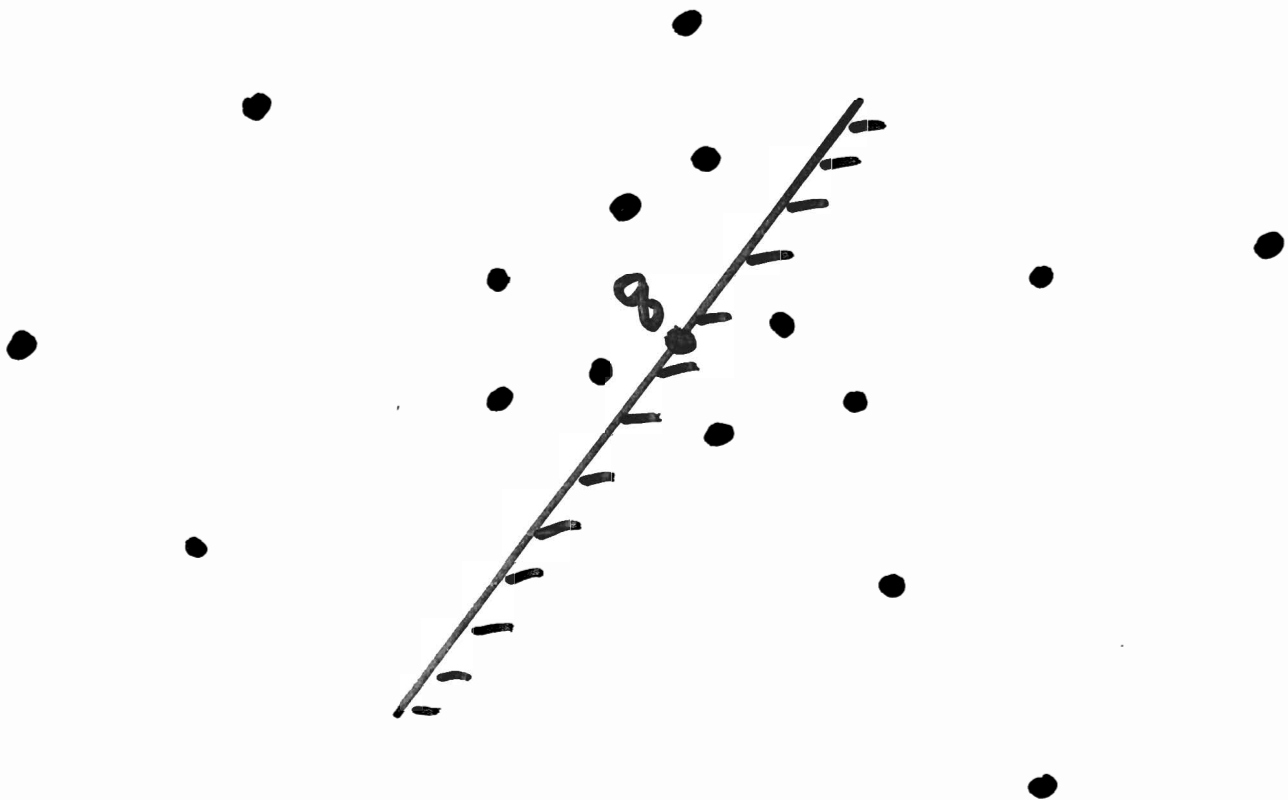
Fact: If R_j is α -sketch of S_j
($|R_1| = |R_2|$, $|S_1| = |S_2|$)
then δ -sketch of $R_1 \cup R_2$ is
 $(\alpha + \delta)$ -sketch of $S_1 \cup S_2$

Set $\delta = \epsilon / \log n$

\Rightarrow space $O\left(\frac{1}{\delta} \log n\right) = \underline{\underline{O\left(\frac{1}{\epsilon} \log^2 n\right)}}$

Example 1: Approx Centerpoint

Given n pts in \mathbb{R}^d ,
find a pt $q \in \mathbb{R}^d$ st.
any halfspace containing q
contains $\geq \frac{n}{d+1} - \epsilon n$ pts



Bagchi, Chaudhary, Eppstein, Goodrich's Alg'm ('04)

Idea: same method!

Def: [Vapnik, Chervonenkis '71 ... Matoušek '95]

RCS is a δ -approximation of S iff
any halfspace containing i pts of R
contains $\sim \frac{i n}{r} \pm \delta n$ pts of S

Fact: \exists δ -approximation of size
 $O\left(\frac{1}{\delta^2} \log \frac{1}{\delta}\right)$

\Rightarrow space $O\left(\frac{1}{\epsilon^2} \text{poly} \log n\right)$

Other Applications [BCEG'04]

more statistics problem

(simplicial depth, LMS, ...)

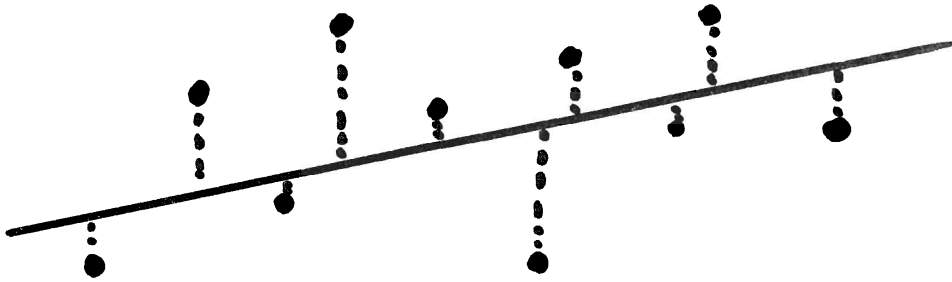
range counting

Example 2: Approx Hyperplane Fitting

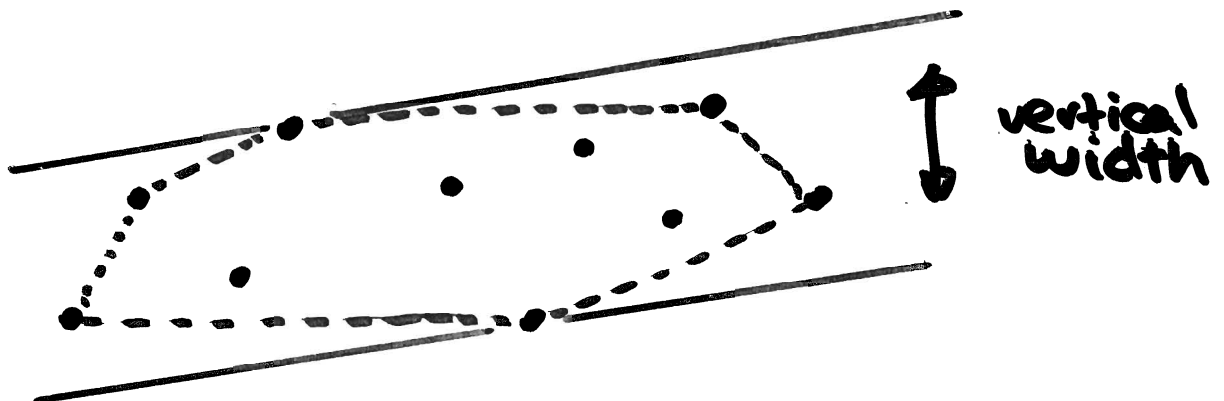
Given n pts in \mathbb{R}^d ,

find hyperplane s.t.

max vertical distance to the pts
is minimized, to within $(1 \pm \epsilon)$ factor



i.e. minimize vertical-width of
convex hull (CH) over all directions



Agarwal, Har-Peled, Varadarajan's Alg'm [04]

Idea: logarithmic method again!

Def: RCS is a δ -core-set of S iff

\forall direction x

width of $CH(S)$ along x

$\leq (1 + \delta) \cdot$ width of $CH(R)$ along x

Fact: \exists δ -core-set of size $O\left(\frac{1}{\delta^{(d-1)/2}}\right)$

\Rightarrow space $O\left(\frac{1}{\epsilon^{O(d)}} \text{poly } \log n\right)$

A New Alg'm [C'04]

Idea: "geometric-series" method

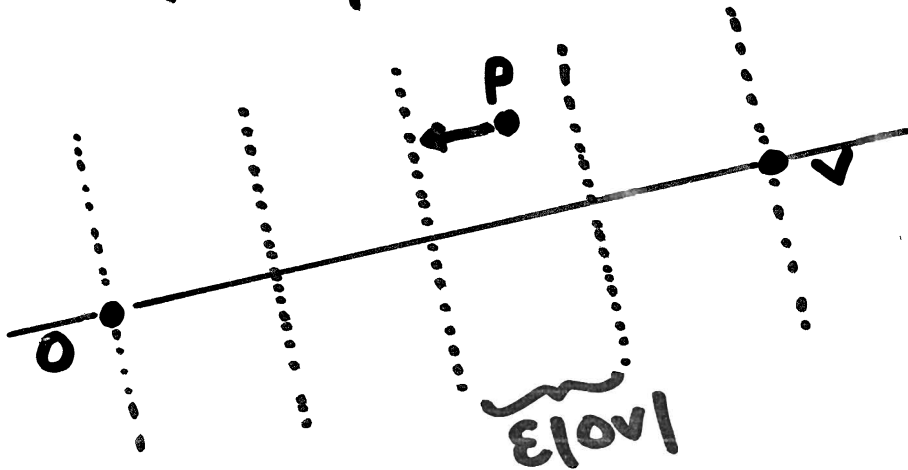
Static version: // core-sets in 2D

o = first pt

v = (approx) farthest pt from o

round pts to $\Theta(1/\epsilon)$ grid lines
orthogonal to \overrightarrow{ov}

keep min/max pts on each line



Analysis: rounding error

$\leq \epsilon \cdot \text{width of } \overline{ov} \text{ along } x$

$\leq \epsilon \cdot \text{width of } CH(S) \text{ along } x$

Modified streaming version:

To insert p :

whenever $|op| > 2|ov|$,

$$v = p$$

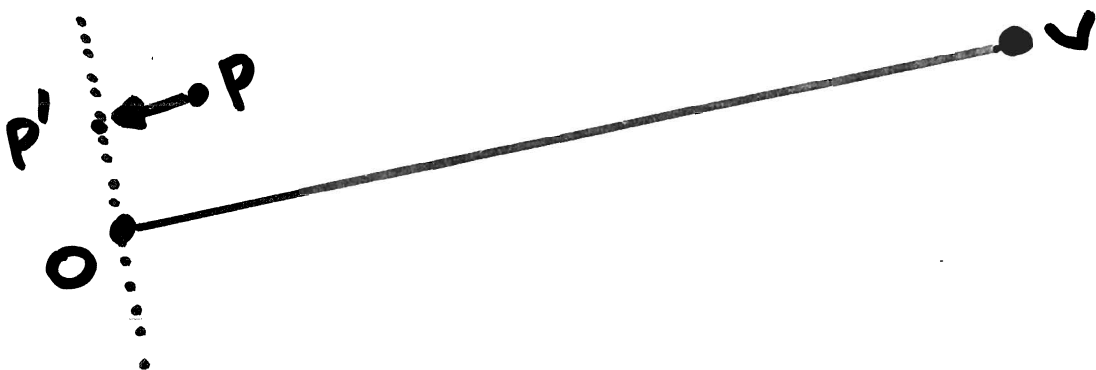
start a new core-set & clean up

To clean up:

keep only the b newest core-sets

round old pts to grid line thru o

keep min/max pt on line



Analysis:

$$\text{rounding error} \leq \frac{|pp'|}{|ov|} \cdot \text{width of } \overline{ov} \text{ along } x$$

$$\leq \frac{|op|}{|ov|} \cdot \text{ " " " }$$

when p was inserted: $|op| \leq 2|ov|$

at every change: $|ov|$ doubles

when p gets old: $|op| \leq \frac{2|ov|}{2^b}$

$$\Rightarrow \text{total error} \leq 2 \left(\frac{1}{2^b} + \frac{1}{2^{b+1}} + \dots \right) \cdot \text{width of } CH(S) \text{ along } x$$

$$\text{Set } b = \log \frac{1}{\epsilon}$$

$$\Rightarrow \text{space } O\left(b \cdot \frac{1}{\epsilon}\right) = \underline{\underline{O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)}} \\ \text{const (indep of } n \text{)!}$$

Other Results & Applications [c'04]

refinement in 2D:

$$\text{space } O\left(\frac{1}{\sqrt{\epsilon}} \log^2 \frac{1}{\epsilon}\right)$$

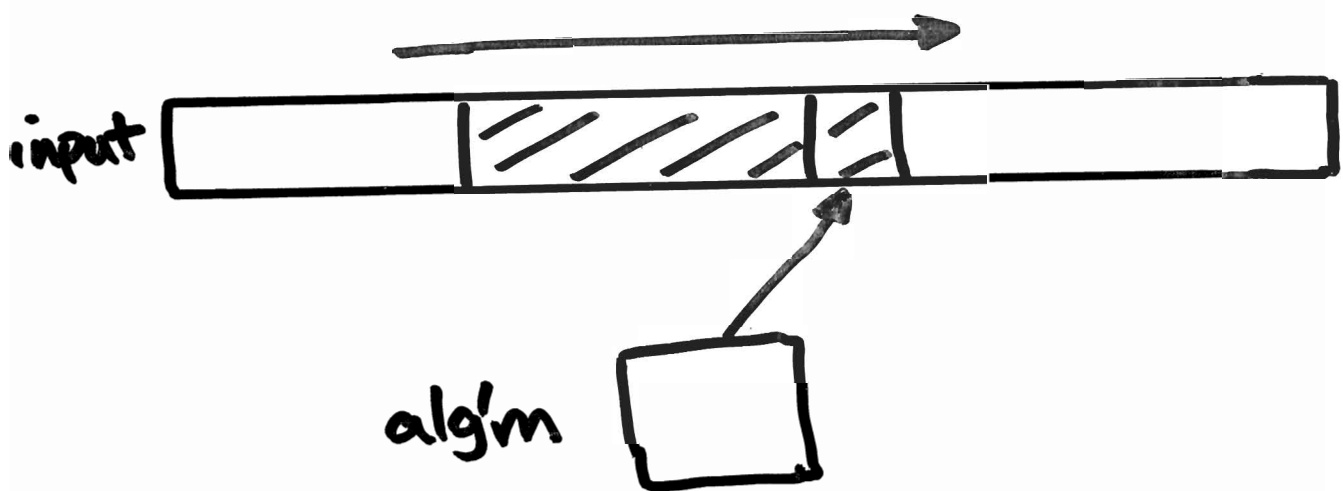
core-sets in dD :

$$\text{space } O\left(\frac{1}{\epsilon^{d-0.5}} \log^d \frac{1}{\epsilon}\right)$$

bounding box

sphere/cylinder fitting

SLIDING-WINDOW ALGMS

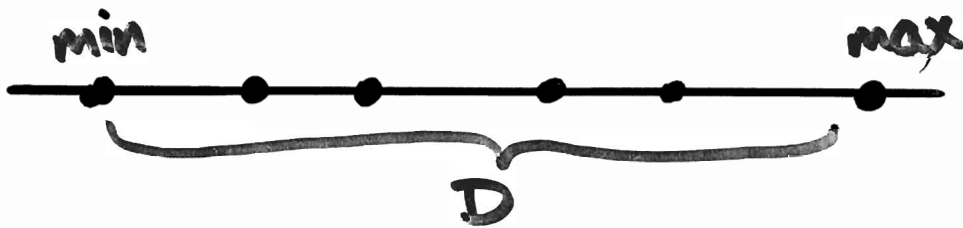


(data structures supporting
"insert" & "delete oldest pt")

Example : Approx Diameter in 1D

Given n pts in \mathbb{R}^1 ,

compute $D =$ farthest distance
to within 1ϵ factor



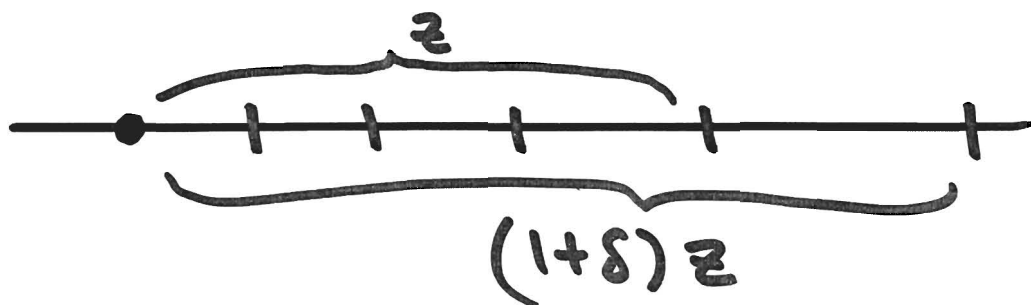
Feigenbaum, Kannan, Zhang's Algm [02]

Idea: logarithmic method again!

How to sketch:

forms exponentially-spaced grid

keeps newest pt (the representative)
of each cell



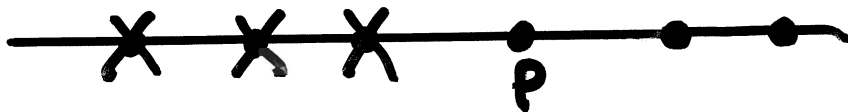
$$\Rightarrow \text{space } \underline{\underline{O\left(\frac{1}{\epsilon} \log \Delta \log^2 n\right)}}$$

where $\Delta = \frac{\text{farthest dist}}{\text{closest dist}}$

A New (Simple!) Alg'm [C, Sadjad '04]

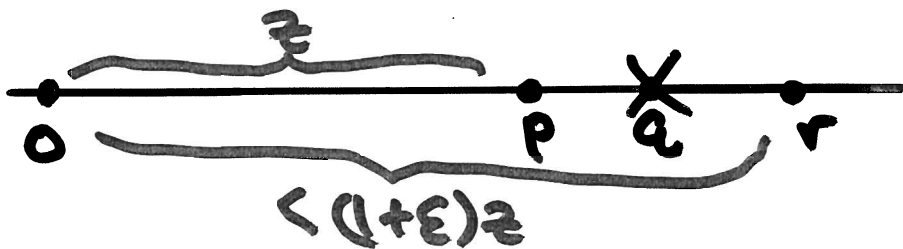
Idea: "skip-middle" method

To insert p : // approx max
insert p to R
remove all pts $< p$ from R
after $\frac{1}{\epsilon} \log \Delta$ inserts, clean up



To clean up:

o = smallest pt of R
whenever \exists 3 consecutive $p, q, r \in R$
s.t. $|or| < (1+\epsilon) |op|$,
delete q from R

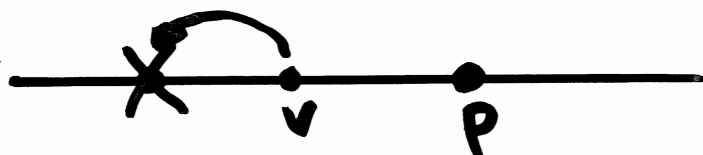


\Rightarrow space $|R| = \underline{\underline{O\left(\frac{1}{\epsilon} \log \Delta\right)}}$ optimal!

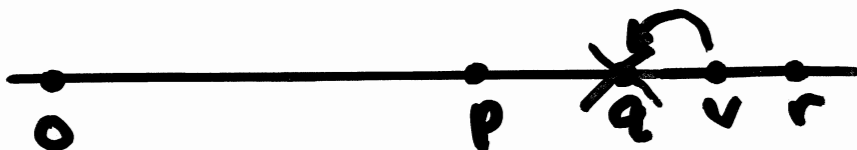
Correctness/Invariants:

1. smaller pts are newer in R
2. each pt v has a representative $v' \in R$ that is newer than v , where either
 - 2.1. v' is successor of v in R , or
 - 2.2. v' is predecessor " " " " with $|vv'| \leq \epsilon \cdot \text{diameter of all pts newer than } v$

(insert)



(clean-up)



$$(|or| \leq (1+\epsilon) |op| \Rightarrow |pv| \leq \epsilon |op|)$$

Other Applications [CS'04]

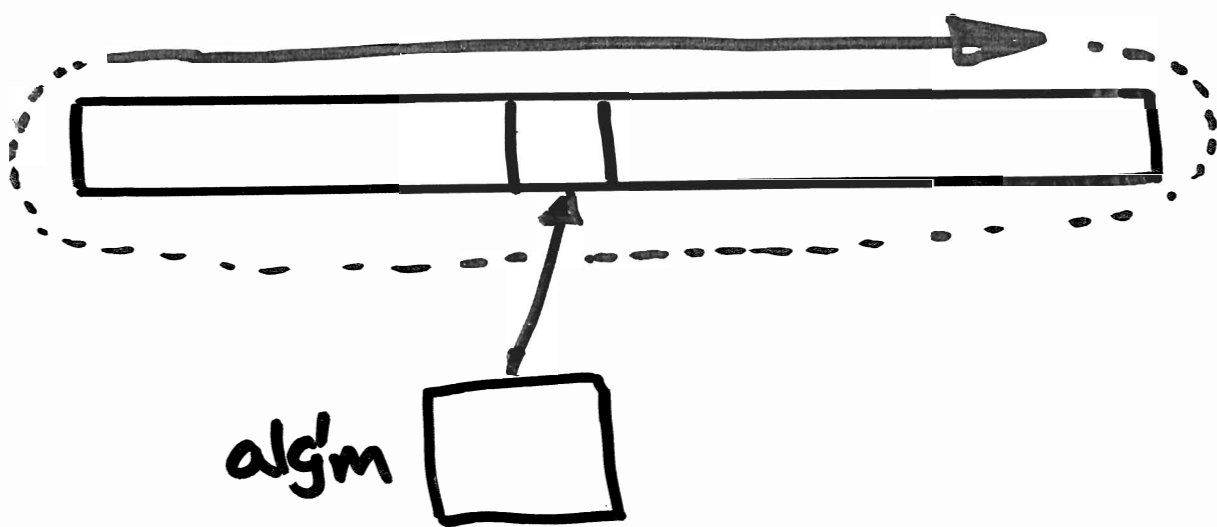
diameter in dD :

$$\text{space } O\left(\frac{1}{\epsilon^{(d-1)/2}} \log \frac{\Delta}{\epsilon}\right)$$

[by projection]

core-sets in 2D

MULTI-PASS ALG'MS



1 Pass vs. 2 Passes,
according to Knuth [AOCF, vol 1, '73]

Old lady, on a bus:

"Can you tell me how I can get
off Pasadena St.?"

Boy:

"Just watch me, and get off
two stops before I do."

Example 0 : Exact (!) Median



Munro, Paterson's Algm [180]

Idea: filtering ("prune-and-search")

$$I = (-\infty, \infty)$$

repeat

1. among pts inside I ,
compute $(1/r)$ -sketch R of size $O(r)$
2. $I =$ sub-interval containing answer

$$\text{Set } r = n^\delta$$

$$\Rightarrow \text{passes } \log_r n = \frac{1}{\delta} \text{ const!}$$

$$\text{space } O(r \log^2 n) = \tilde{O}(n^\delta) \text{ small!}$$

Example 2 : Exact Hyperplane Fitting



$$\begin{aligned} \min t \\ \text{s.t. } y_i &\leq mx_i + b + t \\ y_i &\geq mx_i + b - t \\ (i = 1, \dots, n) \end{aligned}$$

reduces to linear programming (LP)
in \mathbb{R}^{d+1}

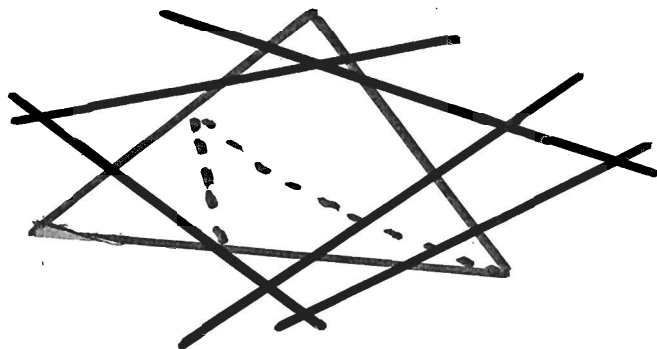
"New" Alg'm [C, Chen '05]

Idea: prune-&-search again

$\Delta = \mathbb{R}^d$ // LP in dD

repeat:

1. among halfspaces crossing Δ , compute $(1/r)$ -cutting of size $r^{O(1)}$
2. $\Delta =$ sub-simplex containing answer



Def: Given n halfspaces, δ -cutting is a partition of space into simplices each crossed by $\leq \delta n$ halfspaces

Analysis:

1. take $(1/r)$ -approximation
2. solve $r^{O(1)}$ LPs in $(d-1)D$
"in parallel"

passes $P_d(n) = O(P_{d-1}(n) \log_r n)$

space $S_d(n) = O(S_{d-1}(n) r^{O(1)} \log_r n)$

Set $r = n^{c\delta}$

\Rightarrow passes $\underline{O(1/\delta^{d-1})}$ const!

space $\underline{\tilde{O}(n^\delta)}$ small!

Other Results [cc'05]

refinement: $O(n)$ rand time

$O(n \log \log \dots n)$ det time in 2D

lower bd for LP in 2D:

$\frac{1}{\delta}$ passes $\Rightarrow \Omega(n^\delta)$ space
(as in Munro, Paterson)

CH of sorted pts in 2D:

$O(1)$ passes, $O(n^{1/2+\delta})$ space,

$O(n)$ time

CH for small output size h in 2D:

$O(1)$ passes, $O(h n^\delta)$ space,

$O(n \log n)$ time

CONCLUSION

Summary of New Results

one-pass alg's for core-sets
(approx CH) [C'04]

sliding-window alg's for
diameter [CS'04]

multi-pass alg's for LP [CC'05]

Some Open Problems

one-pass for high dim

e.g. diameter: $\sim \sqrt{2}$ factor [Indyk'03]

min enclos cylinder:

~ 5 factor [C'04]

smallest # passes

e.g. LP: $d+1$ passes, $\tilde{O}(\sqrt{n})$ space
[Clarkson's alg'm]

lower bd?