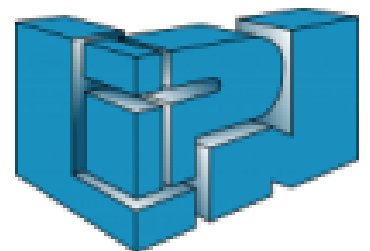


Small Random Samples via the Crossing Distance

Metric Graph Theory and Related Topics

NABIL H. MUSTAFA

LIPN, USPN



UNIFORM APPROXIMATIONS OF SET SYSTEMS

APPROXIMATIONS OF SET SYSTEMS

n elements m sets
 ↙ ↘
 set system (X, \mathcal{R})

compute an approximation $A \subseteq X$

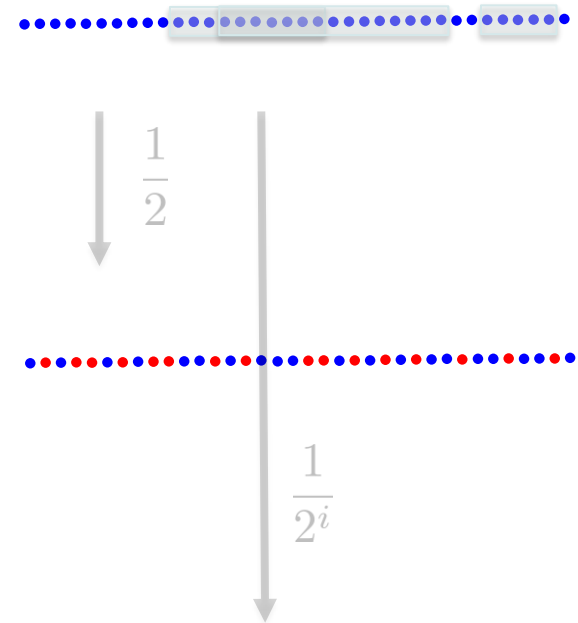
$S \in \mathcal{R}$

$|A|$ $\mathbf{E}[|A \cap S|]$

$\frac{n}{2}$ $\frac{|S|}{2}$ $\frac{|S|}{2} \pm \boxed{\text{error}}$ discrepancy

i steps

$\boxed{\frac{n}{2^i}}$ $\frac{|S|}{2^i}$ $\frac{|S| |A|}{n} \pm \epsilon |A|$ ϵ -approximation



RANDOM SAMPLING FOR DISCREPANCY

n m
set system (X, \mathcal{R})

Chernoff-Hoeffding bound:

For any $\Delta > 0$ and $S \in \mathcal{R}$

$$\Pr[Y_S \geq \Delta] < \exp\left(-\frac{\Delta^2}{2|S|}\right)$$

color each point $\{+1, -1\}$ independently at random

for a fixed $S \in \mathcal{R}$: discrepancy $O\left(\sqrt{|S|}\right)$ with constant probability

probability of discrepancy
being at least η for some $S \in \mathcal{R}$

$$< m \cdot \exp\left(-\Theta\left(\frac{\eta^2}{n}\right)\right)$$
$$\leq m \cdot e^{-\ln 2m} \leq \frac{1}{2}$$

discrepancy \rightarrow $\eta = \Theta\left(\sqrt{n \ln m}\right)$

RANDOM SAMPLING FOR APPROXIMATIONS

set system (X, \mathcal{R}) n m **discrepancy** $O\left(\sqrt{n \ln m}\right)$

for all $S \in \mathcal{R}$:

$$A_1 \subseteq X$$

$$|A_1 \cap S| = \frac{|S|}{2} \pm O\left(\sqrt{n \ln m}\right)$$

$$A_2 \subseteq A_1$$

$$|A_2 \cap S| = \frac{|A_1 \cap S|}{2} \pm O\left(\sqrt{\frac{n}{2} \ln m}\right) = \frac{|S|}{4} \pm O\left(\frac{1}{2}\sqrt{n \ln m} + \sqrt{\frac{n}{2} \ln m}\right)$$

after i steps:

$$|A_i \cap S| = \frac{|S|}{2^i} \pm O\left(\sqrt{\frac{n}{2^i} \ln m}\right) \leq \epsilon \frac{n}{2^i}$$

setting $t = \log \frac{\epsilon^2 n}{\ln m}$ gives

$$|A_t| = O\left(\frac{1}{\epsilon^2} \log m\right)$$

$$|A_t \cap S| = \frac{|S| |A_t|}{n} \pm \epsilon |A_t|$$

LOCALLY NICE SYSTEMS

Theorem: A uniform random sample of X of size $\Theta\left(\frac{1}{\epsilon^2} \ln m\right)$ is an ϵ -approximation.

‘locally polynomial’ set systems

total number of sets $|\mathcal{R}|$: $O(n^4)$

number of subsets on $Y \subseteq X$: $O(n^4) \rightarrow O(|Y|^4)$

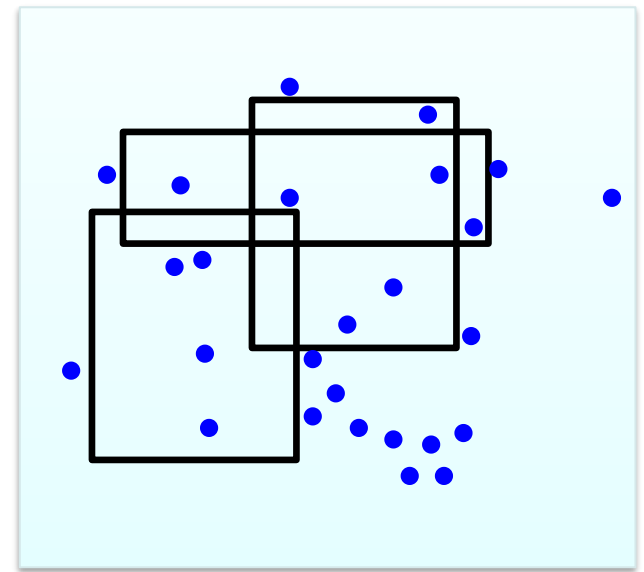
combinatorially

$$\mathcal{R}|_Y = \{Y \cap R : R \in \mathcal{R}\}$$

the *projection* of \mathcal{R} onto Y

a constant d such that

$$|\mathcal{R}|_Y| = O(|Y|^d) \text{ for any } Y \subseteq X$$



$X \quad \mathcal{R}$

RANDOM SAMPLING FOR APPROXIMATIONS

set system (X, \mathcal{R}) n $m = O(n^d)$ **discrepancy** $O(\sqrt{n \ln m})$

for all $S \in \mathcal{R}$:

$$A_1 \subseteq X$$

$$|A_1 \cap S| = \frac{|S|}{2} \pm O\left(\sqrt{n \ln m}\right) n^d$$

$$A_2 \subseteq A_1$$

$$|A_2 \cap S| = \frac{|A_1 \cap S|}{2} \pm O\left(\sqrt{\frac{n}{2} \ln m}\right) = \frac{|S|}{4} \pm O\left(\frac{1}{2} \sqrt{n \ln m} + \sqrt{\frac{n}{2} \ln m}\right) \frac{(n/2)^d}{n^d}$$

after i steps:

$$|A_i \cap S| = \frac{|S|}{2^i} \pm O\left(\sqrt{\frac{n}{2^i} \ln m}\right) (n/2^i)^d$$

setting $t = \log \frac{\epsilon^2 n}{\ln m}$ gives $\left(\frac{1}{\epsilon}\right)^d$

$$|A_t| = O\left(\frac{1}{\epsilon^2} \log m\right) \left(\frac{1}{\epsilon}\right)^d$$

$$|A_t \cap S| = \frac{|S| |A_t|}{n} \pm \epsilon |A_t|$$

APPROXIMATION BOUNDS

n elements

m subsets

d dimension

Random
Arbitrary

VC dimension

Uniform Sampling

Discrepancy

$$\sqrt{n \ln m}$$

$$\sqrt{dn \ln n}$$

Approximations

$$\frac{1}{\epsilon^2} \ln m$$

$$\frac{d}{\epsilon^2} \ln \frac{1}{\epsilon} \rightarrow \frac{d}{\epsilon^2}$$

[Talagrand, 1994]

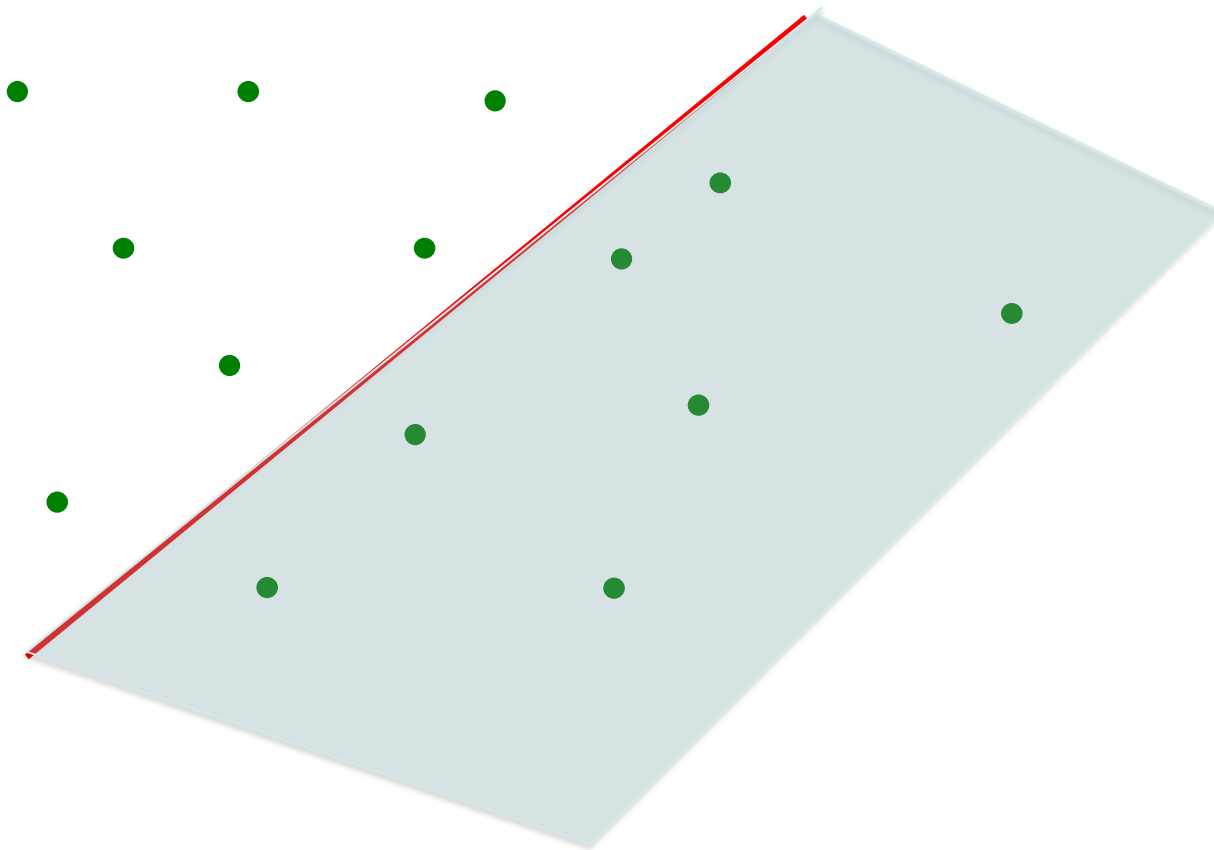
[Li, Long, Srinivasan, 2001]

usefulness: **oblivious uniform approximation**

NON-UNIFORM APPROXIMATIONS

SET SYSTEM INDUCED BY HALF-SPACES

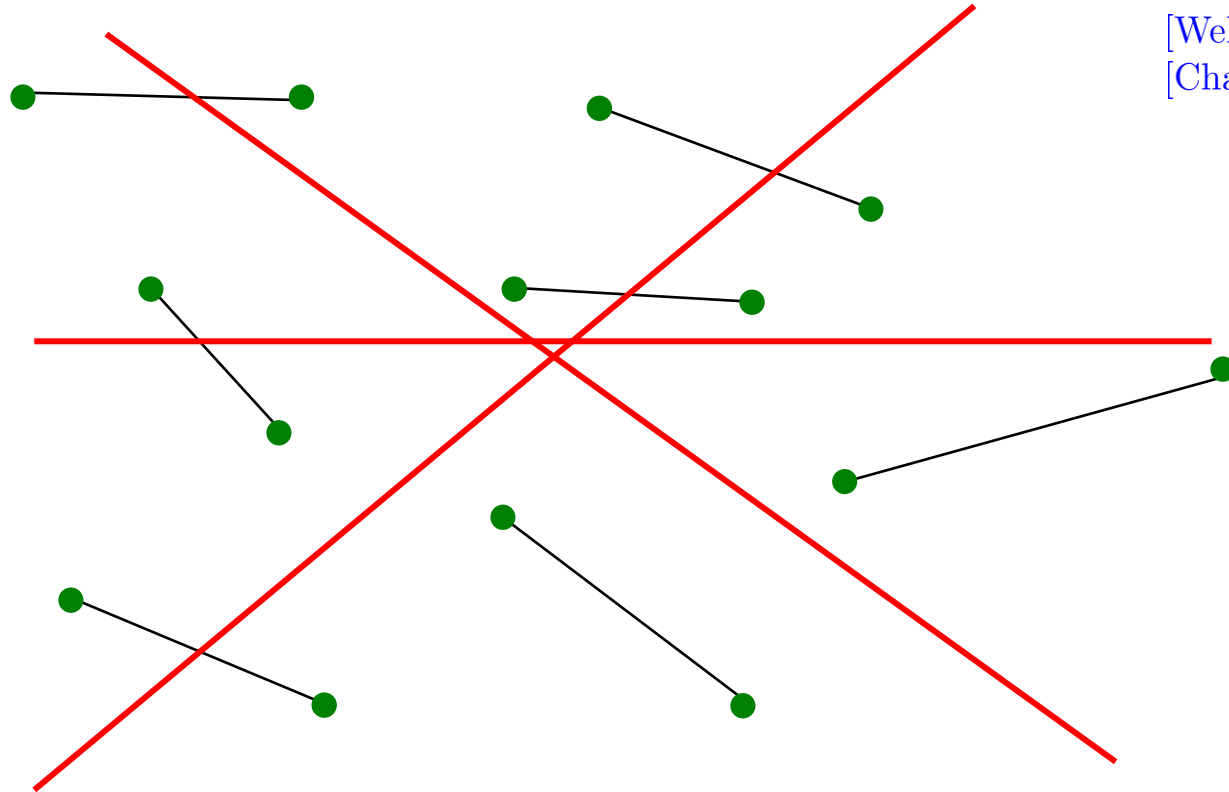
set system induced by half-spaces in \mathbb{R}^d



$$O\left(\sqrt{n \ln n^2}\right) \rightarrow O\left(n^{1/4}\right) \quad \text{optimal}$$

SPANNING TREES WITH LOW CROSSING NUMBER

Theorem: Given any set P of n points in \mathbb{R}^2 , there exists a matching M on P such that any line crosses at most $O(\sqrt{n})$ edges of M .



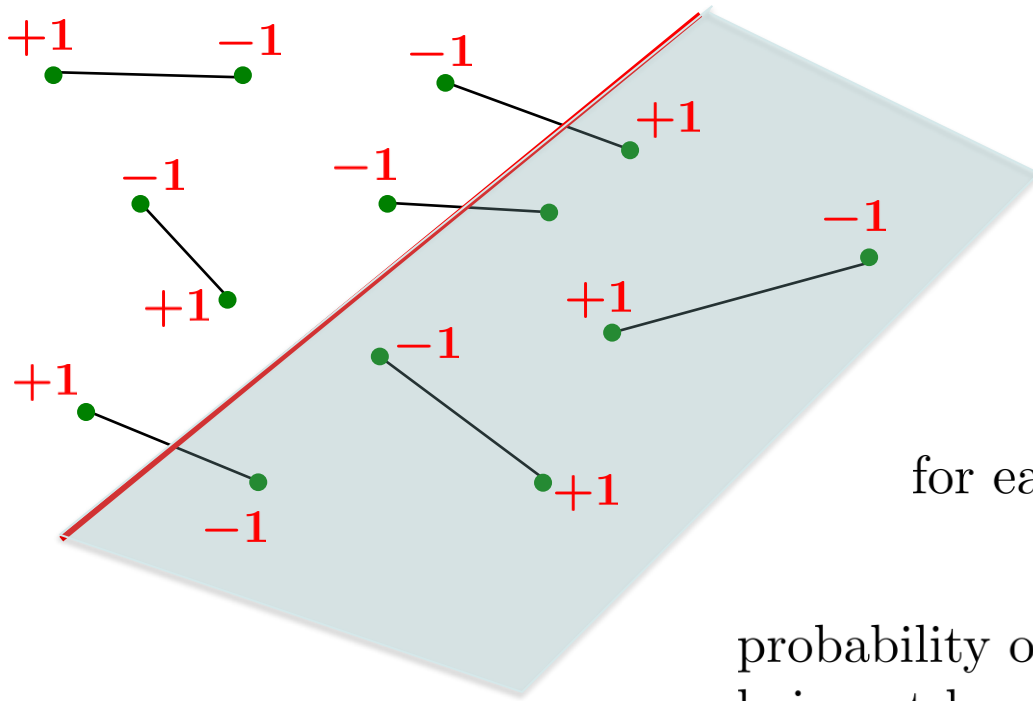
[Welzl 1988]

[Chazelle-Welzl 1989]

One almost wouldn't believe that after thousands of years of geometry, it is still possible to discover such pretty theorems about points in the plane.

J. Matoušek, ICM 1998

SPANNING TREES WITH LOW CROSSING NUMBER



Chernoff-Hoeffding bound:

For any $\Delta > 0$ and $S \in \mathcal{R}$

$$\Pr[Y \geq \Delta] < \exp\left(-\frac{\Delta^2}{2|S|}\right)$$

for each half-space, $O(\sqrt{n})$ events

probability of discrepancy
being at least η

$$\leq m \cdot \exp\left(-\Theta\left(\frac{\eta^2}{\sqrt{n}}\right)\right)$$

$$\leq m \cdot e^{-\ln 2m} \leq \frac{1}{2}$$

$$\eta = \Theta\left(n^{1/4} \sqrt{\ln m}\right)$$

APPROXIMATION BOUNDS

n elements
 m subsets
 d dimension

Random

VC dimension

Uniform Sampling

Non-Uniform Sampling
 +
 Combinatorics

Discrepancy

$$\sqrt{n \ln m}$$

$$\sqrt{dn \ln n}$$

$$n^{\frac{1}{2}} - \frac{1}{2d}$$

Approximations

$$\frac{1}{\epsilon^2} \ln m$$

$$\frac{d}{\epsilon^2}$$

$$\frac{1}{\epsilon^{2 - \frac{2}{d+1}}}$$

Problem: now need to construct matchings to be able to sample

faster constructions imply improved algorithms for constructing small samples

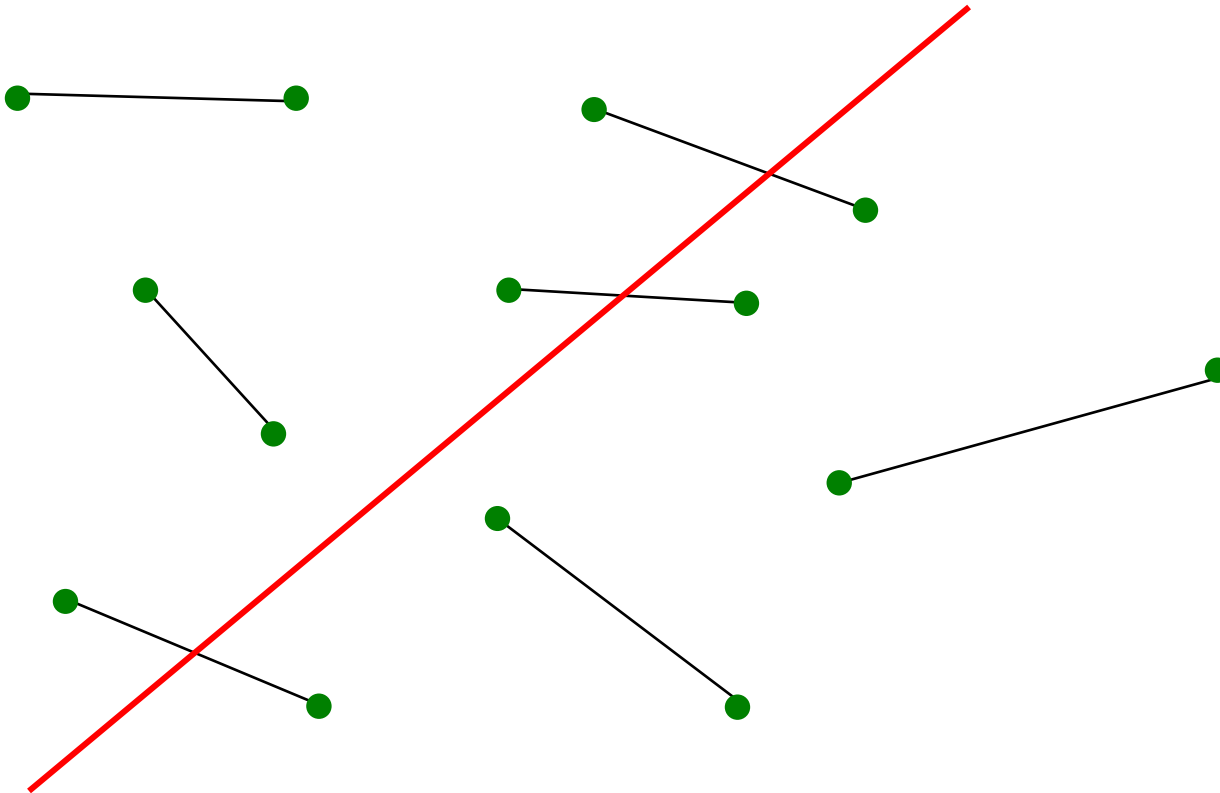
LOW CROSSING MATCHINGS

SPANNING TREES WITH LOW CROSSING NUMBER

Theorem: Given any set P of n points in \mathbb{R}^2 , there exists a matching M on P such that any line crosses at most $O(\sqrt{n})$ edges of M .

[Welzl 1988]

[Chazelle-Welzl 1989]

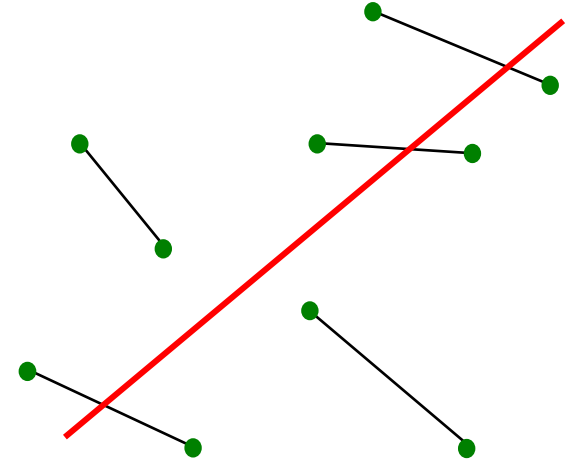


More generally: construct matchings for set systems, with a natural generalized notion of crossing

SPANNING TREES WITH LOW CROSSING NUMBER

Given points P , and lines \mathcal{L} , let M be a tree on P such that any line in \mathcal{L} crosses at most $O(\sqrt{n})$ edges of M .

$$|P| = n \quad |\mathcal{L}| = m$$



each line $l \in \mathcal{L}$ crosses at most \sqrt{n} edges of M

summed over all $l \in \mathcal{L}$, there are $m \cdot \sqrt{n}$ crossings between lines and edges

one edge of M must be crossed by at most $\frac{m \cdot \sqrt{n}}{n-1} = O\left(\frac{m}{\sqrt{n}}\right)$ lines in \mathcal{L}

Key Property: Given P and \mathcal{L} , there exist two points $p_i, p_j \in P$ such that the line segment $\overline{p_i p_j}$ crossed by $O\left(\frac{m}{\sqrt{n}}\right)$ lines in \mathcal{L} .

SPANNING TREES WITH LOW CROSSING NUMBER

Key Property: Given X and \mathcal{S} , there exist two elements $p_i, p_j \in X$ such that the pair $\{p_i, p_j\}$ is crossed by $\leq \frac{m}{\sqrt{n}}$ sets in \mathcal{S} .

$X_p \subseteq \mathcal{S}$: sets of \mathcal{S} containing p

need to show: two sets X_p, X_q with symmetric difference $O\left(\frac{m}{\sqrt{n}}\right)$

if symmetric difference between every pair is $\Omega(\Delta)$

Packing lemma: number of sets is then $O\left(\left(\frac{m}{\Delta}\right)^2\right)$

[Haussler 1995]

[Ezra 2014]

[M. 2016]

$$O\left(\left(\frac{m}{\Delta}\right)^2\right) < n \quad \rightarrow \quad \Delta > \frac{m}{\sqrt{n}}$$

key is structure of the k -distance graph on (X, \mathcal{S})

$k = 1$ [Haussler-Littlestone-Warmuth 1990]

$k = 2$ [Chepoi-Labourel-Ratel 2019]

much better understood for partial cubes

[Chepoi-Knauer-Philibert 2020]

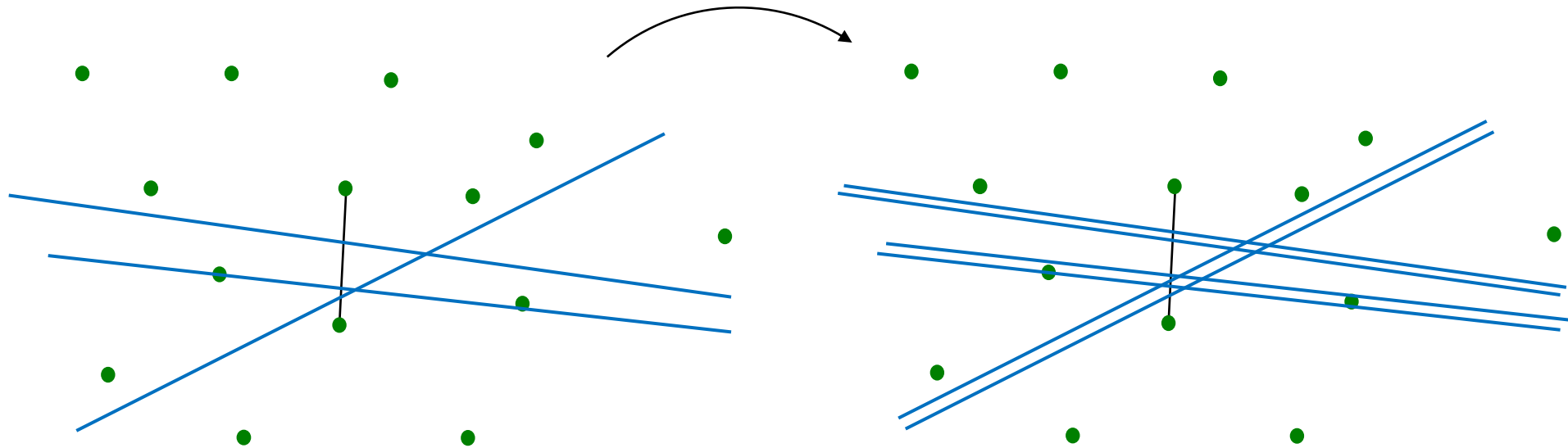
SPANNING TREES WITH LOW CROSSING NUMBER

first idea: construct tree iteratively by adding these edges one by one

given n points and \mathcal{L} , find an edge intersecting $\leq m/\sqrt{n}$ lines in \mathcal{L}

add to M , remove both endpoints, and iterate on $n - 2$ remaining points

second idea: want to discourage a line from intersecting too many edges



at each step, add an edge e , and double the lines that intersect e

iterate on the remaining $n - 2$ points with the new larger set of lines \mathcal{L}'

very surprising that such a heuristic-y argument works

SPANNING TREES WITH LOW CROSSING NUMBER

Why does it work?

at each step, find an edge intersecting at most m/\sqrt{n} lines

so number of new lines added to \mathcal{L} are at most m/\sqrt{n}

at the end, the total number of lines added is ‘small’

Now take any line $l \in \mathcal{L}$

each time it intersects an edge, it’s weight is doubled

if it intersects k edges, total weight in the end: 2^k

but 2^k is upper-bounded by the total weight, which is small

so k cannot be too large!

SPANNING TREES WITH LOW CROSSING NUMBER

[Welzl 1987]

[Chazelle-Welzl 1989]

Points

Lines

Edge Intersects

Step 1 n m m/\sqrt{n}

Step 2 $n - 2$ $m(1 + 1/\sqrt{n})$ $\frac{m(1+1/\sqrt{n})}{\sqrt{n-2}}$

Step 3 $n - 4$ $m(1 + 1/\sqrt{n})(1 + 1/\sqrt{n-2})$ $\frac{m(1+1/\sqrt{n})(1+1/\sqrt{n-2})}{\sqrt{n-4}}$

total weight of lines after $n/2$ steps: $m \prod_{i=0}^{n/2} \left(1 + 1/\sqrt{n-2i}\right)$

if any line intersects k edges, it has weight $2^k \leq m \prod_{i=0}^{n/2} \left(1 + 1/\sqrt{n-2i}\right)$

solving this gives $k = O(\sqrt{n})$

Problem: how to efficiently find the edge to add at each iteration

NEW ALGORITHM

SPANNING TREES WITH LOW CROSSING NUMBER

[Erdős-Selfridge 1973]

[Grigoriadis-Kachiyan 1995]

[Koufogiannakis-Young 2014]

[Agarwal-Pan 2016]

Algorithm:

set $\omega_1(e) = 1$ for all $e \in E$

set $\pi_1(S) = 1$ for all $S \in \mathcal{S}$

For $i = 1, \dots, \frac{n}{4}$:

e_i : sampled from distribution induced by ω_i $\Pr[e_i = e] = \frac{\omega_i(e)}{\omega_i(E)}$

S_i : sampled from distribution induced by π_i $\Pr[S_i = S] = \frac{\pi_i(S)}{\pi_i(\mathcal{S})}$

double weight of each $S \in \mathcal{S}$ crossing e_i

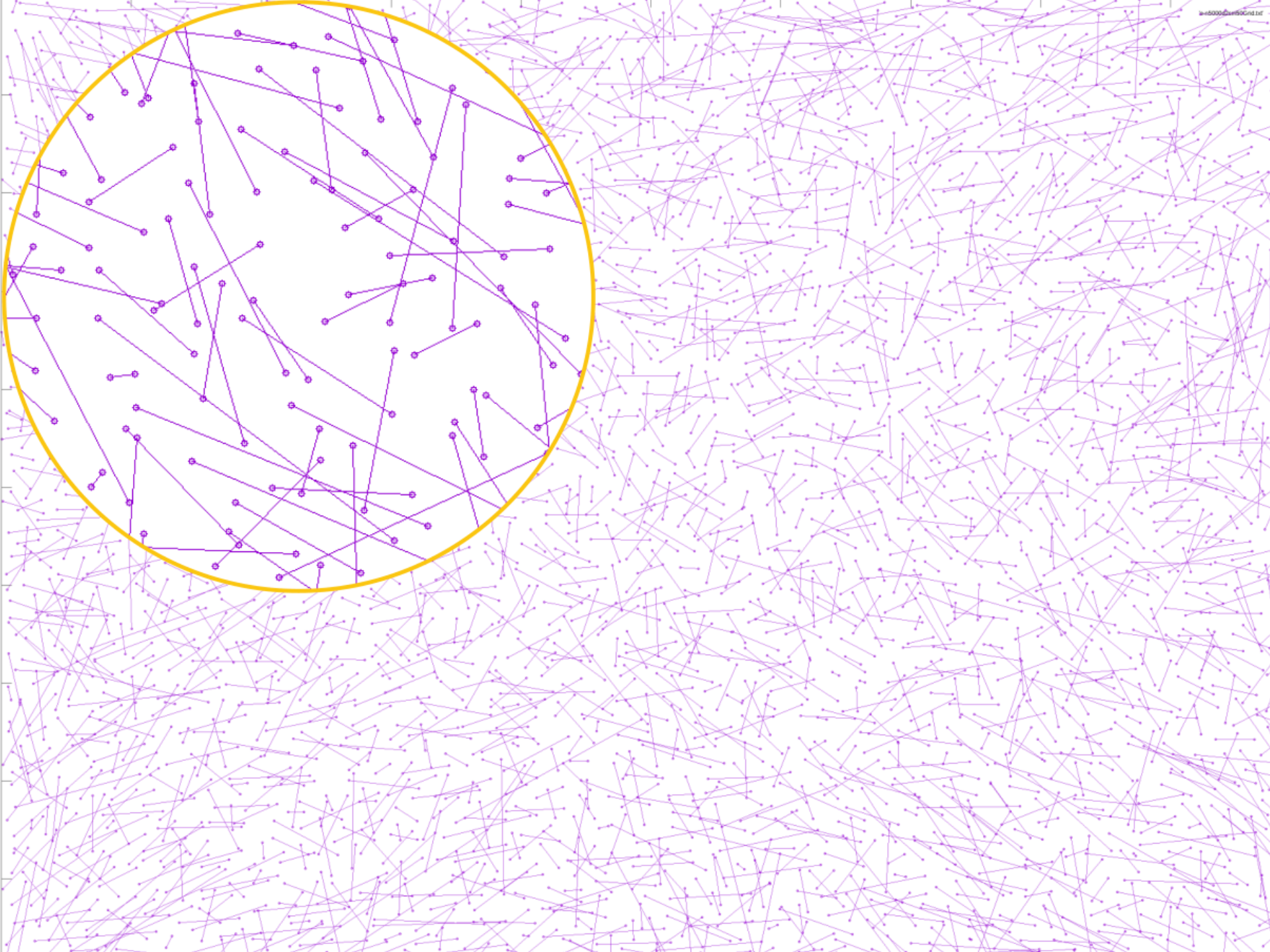
for $S \in \mathcal{S}$: $\pi_{i+1}(S) = \pi_i(S) (1 + I(e_i, S))$

halve weight of each $e \in E$ crossing S_i

for $e \in E$: $\omega_{i+1}(e) = \omega_i(e) (1 - \frac{1}{2}I(e, S_i))$

add e_i to matching

set weights of all edges incident to the two endpoints of e_i to 0



SPANNING TREES WITH LOW CROSSING NUMBER

running time too slow!

[Csikos-M. 2021]

above algorithm improved using structural properties of crossing distance:

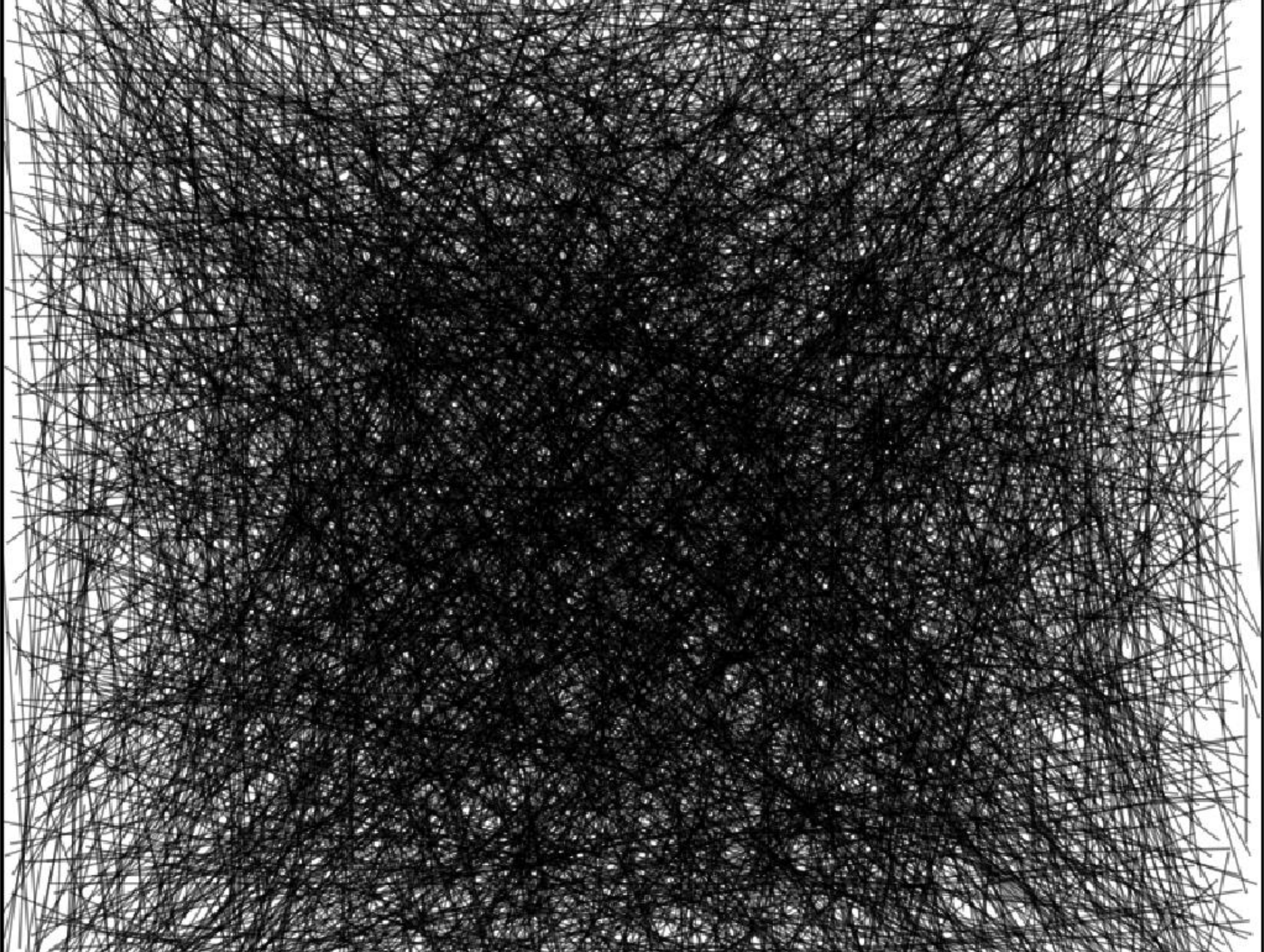
pigeonholing gave a bound on ‘short edge’

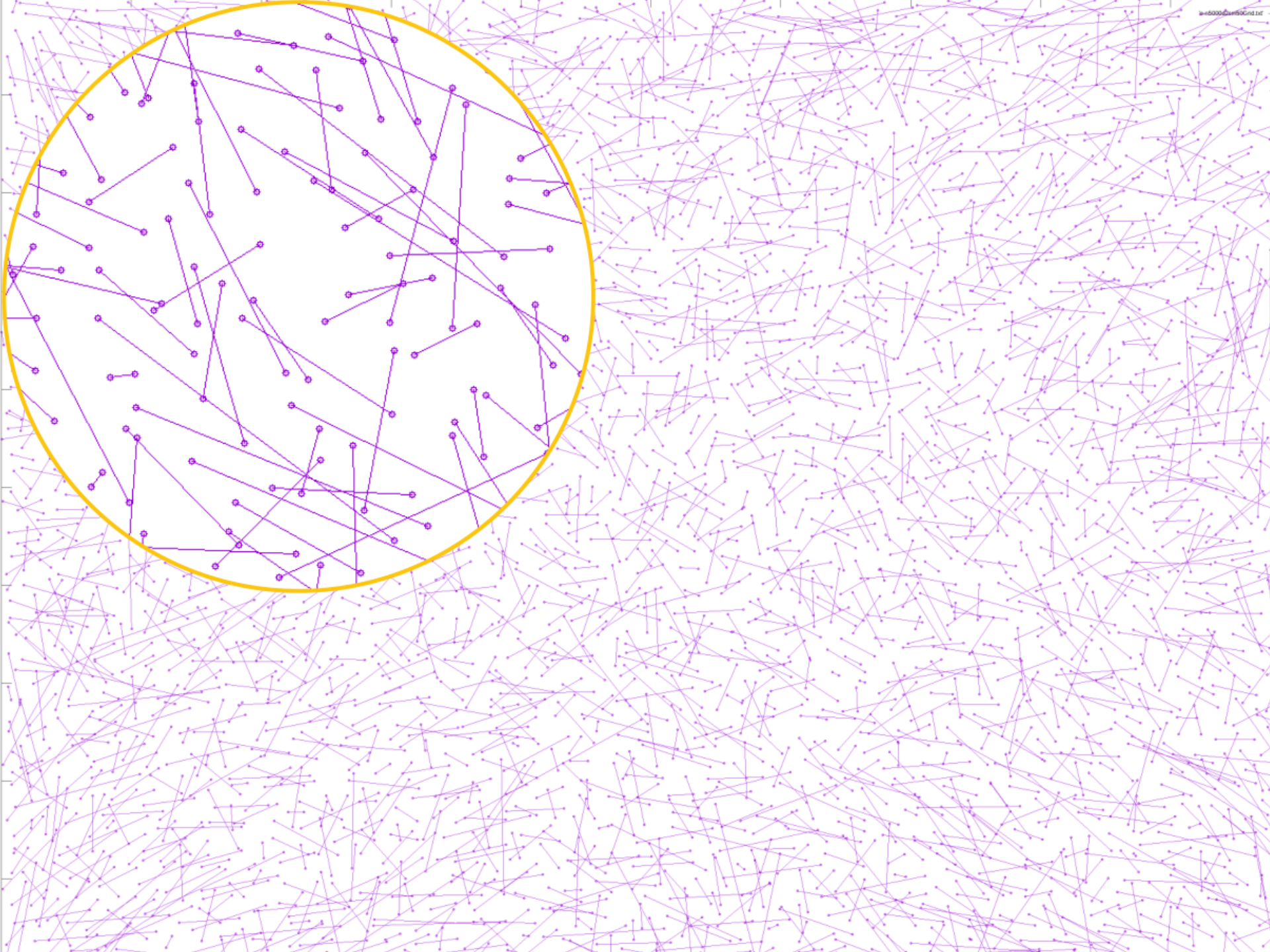
in fact, *many* short edges \rightarrow existence of a sparse set of edges

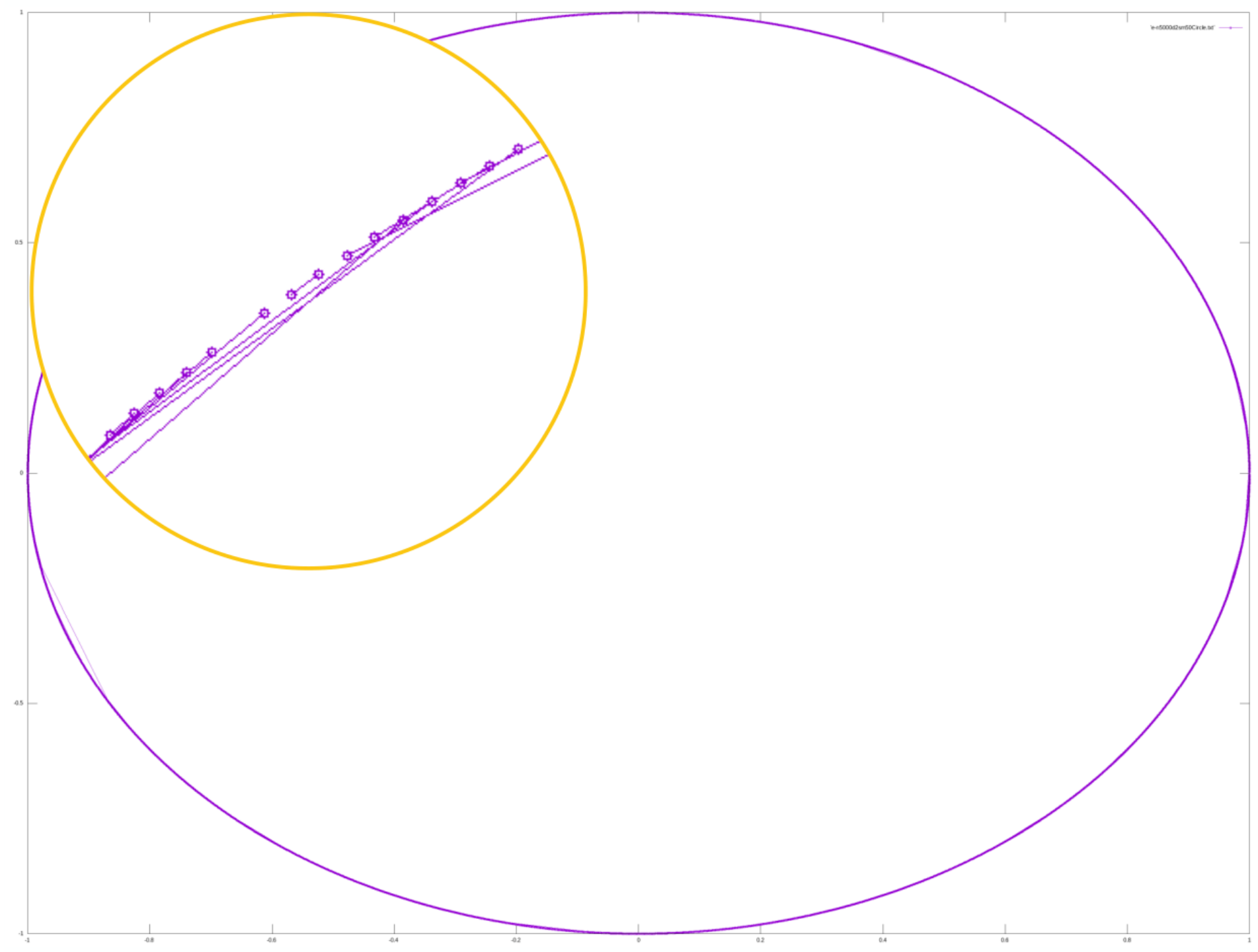
also, able to show that update only a random sample at each step

more involved probabilistic analysis using martingales

simple, easy implementation







[Csikos-M. 2021]

	Our crossing number	Our time	Previous crossing number	Previous time
Arbitrary with $\pi_{\mathcal{S}}^*(k) \leq ck^d$	$\left(\frac{5c^{1/d}d}{d-1} + o(1) \right) n^{1-\frac{1}{d}}$	$\tilde{O} \left(mn^{2/d} + n^{2+2/d} \right)$	$O \left(n^{1-1/d} \right)$	$\tilde{O} \left(mn^2 \right)$ [Har-Peled 2009] [Chekuri-Quanrud 2018]
Geometric induced by balls in \mathbb{R}^d	$\left(6d^2 + o(d^2) \right) n^{1-\frac{1}{d}}$	$\tilde{O} \left(dn^{2+2/d} \right)$	$O \left(n^{1-1/d} \right)$	$\tilde{O} \left(n^{3+1/d} \right)$ [Har-Peled 2009] [Chekuri-Quanrud 2018]
Geometric induced by semi-alg. ranges in \mathbb{R}^d (s eq.'s of deg Δ)	$\left(\frac{20e\Delta sd}{d-1} + o(1) \right) n^{1-\frac{1}{d}}$	$\tilde{O} \left(s\Delta^d \left(mn^{2/d} + n^{2+2/d} \right) \right)$	$O \left(10^d s\Delta n^{1-1/d} \right)$	$O \left(n^{O(d^3)} \right)$ [Agarwal-Matousek-Sharir 2013]
Geometric induced by half-spaces in \mathbb{R}^d	$\left(6d^2 + o(d^2) \right) n^{1-\frac{1}{d}}$	$\tilde{O} \left(dn^{2+2/d} \right)$	$\geq 264d^4 \cdot n^{1-1/d}$	$O \left(d^2 n \right)$ [Chan 2012]

Thank you

**Sampling in Combinatorial
and Geometric Set Systems**

Nabil H. Mustafa

 **AMS** AMERICAN
MATHEMATICAL
SOCIETY