

Finding Haystacks (and Similar Structures) in Geometry

Sariel Har-Peled

²UIUC, Illinois, USA

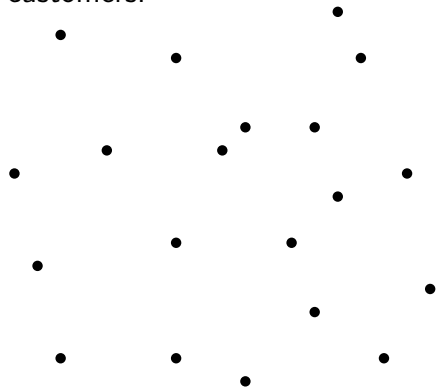
September 2, 2010

2: Motivating problem

Placing an antenna.

P: Set of n points (customers)

Q: Find location of antenna that serves maximum number of customers.

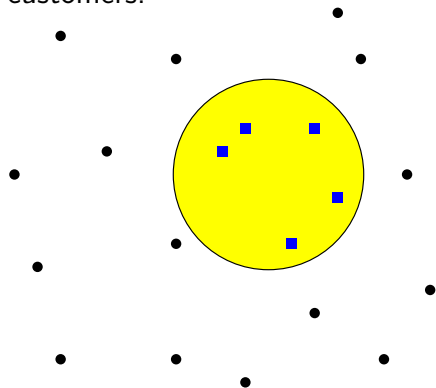


2: Motivating problem

Placing an antenna.

P: Set of n points (customers)

Q: Find location of antenna that serves maximum number of customers.

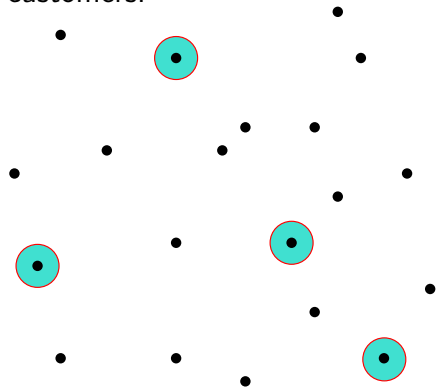


2: Motivating problem

Placing an antenna.

P: Set of n points (customers)

Q: Find location of antenna that serves maximum number of customers.



Approach...

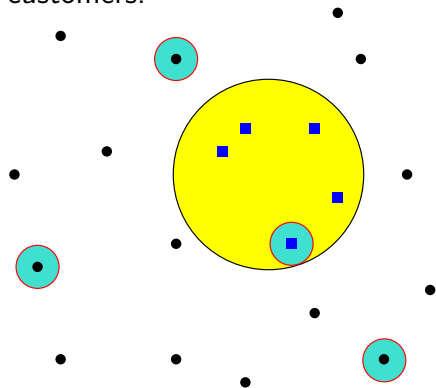
- Pick a random sample.
- Count sample points.

2: Motivating problem

Placing an antenna.

P: Set of n points (customers)

Q: Find location of antenna that serves maximum number of customers.



Approach...

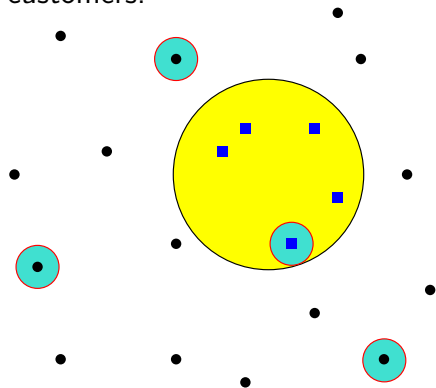
- Pick a random sample.
- Count sample points.

2: Motivating problem

Placing an antenna.

P: Set of n points (customers)

Q: Find location of antenna that serves maximum number of customers.



Approach...

- Pick a random sample.
- Count sample points.

Challenge...

- Good for all disks.
- ... Infinite number of disks.

3: What we want

Definitions

P: Set of points in the plane.

Measure

For any disk **D** its **measure** is $\bar{P}(D) = \frac{|D \cap P|}{|P|}$.

ϵ -Sample

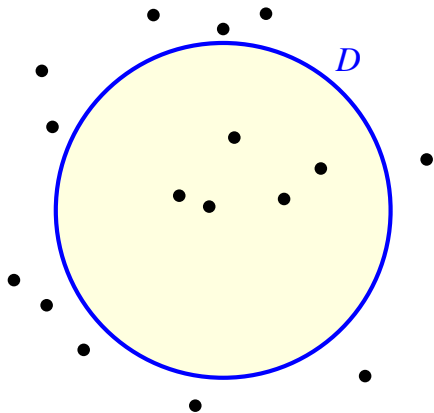
$\epsilon > 0$: Approx parameter.

S \subseteq **P** is ϵ -sample if

$$\forall D \quad |\bar{P}(D) - \bar{S}(D)| < \epsilon$$

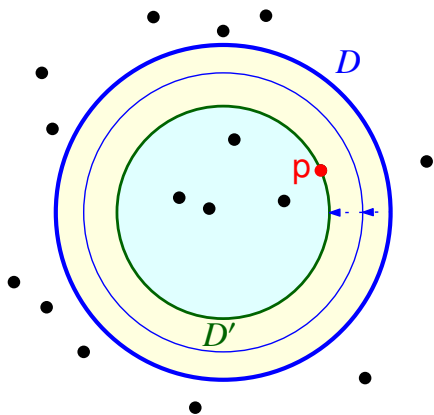
4: Bounding infinity...

...or, it just looks infinite.



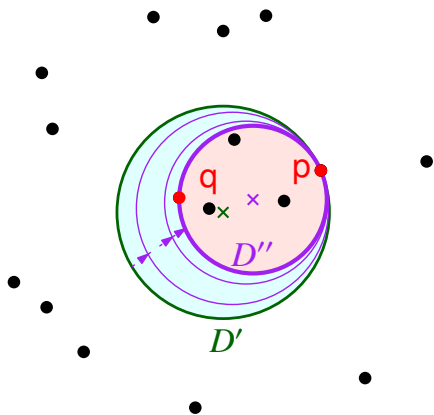
4: Bounding infinity...

...or, it just looks infinite.



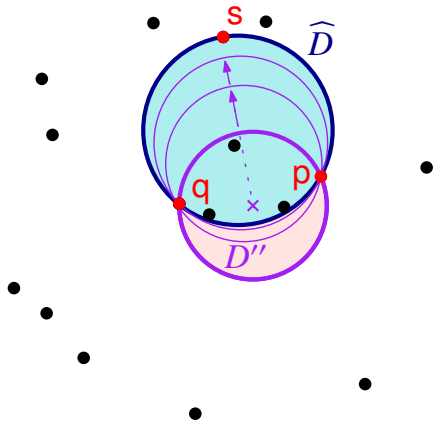
4: Bounding infinity...

...or, it just looks infinite.



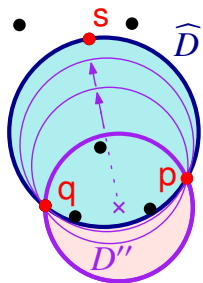
4: Bounding infinity...

...or, it just looks infinite.



4: Bounding infinity...

...or, it just looks infinite.



As such...

n: number of points.

Only $O(n^3)$ different subsets.

Conclusion: ϵ -Sample

Std. random sampling...

$O\left(\frac{\log n}{\epsilon^2}\right)$.

Known: ϵ -Sample

[Vapnik and Chervonenkis, 1971]

[Li et al., 2001] $O\left(\frac{1}{\epsilon^2}\right)$.

5: Application

...or why size matters.

Problem

P: n points in the plane

Compute: Smallest disk containing half the points.

min_disk(**P**, $n/2$): $O(n^2)$ time exact algorithm.

Approximation algorithm

$\epsilon > 0$: Approx parameter.

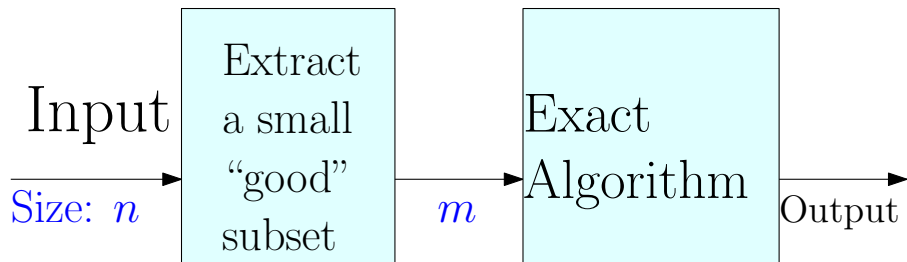
S \subseteq **P**: ϵ -Sample of size m .

Return **min_disk**(**S**, $(1 - \epsilon)m/2$).

Ret disk has $\geq (1/2 - 2\epsilon)n$ points and smaller than exact answer.

6: Application

...or why size matters.



Exact Running time: $T_{\text{exact}}(n)$

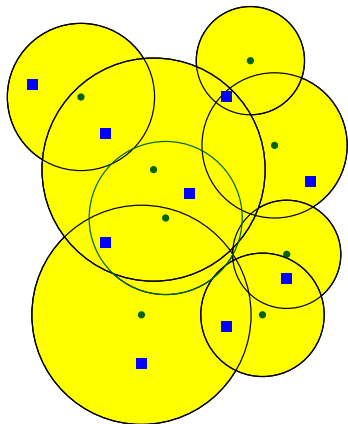
New : $T_{\text{extract}}(n) + T_{\text{exact}}(m)$

Example: Min disk containing $n/2$ points

$T_{\text{extract}}(n) = O(m)$, $T_{\text{exact}}(n) = O(n^2)$, $m = 1/\epsilon^2$.

Running time: $O(1/\epsilon^4)$.

7: Geometric set systems



- \mathcal{F} : m of disks in the plane.
- \mathbf{P} : set of points
- Every disk $\mathbf{D} \in \mathcal{F}$ corresponds to the subset $\mathbf{P} \cap \mathbf{D}$.
- Induced **set system**:

$$\left(\mathbf{P}, \left\{ \mathbf{P} \cap \mathbf{D} \mid \mathbf{D} \in \mathcal{F} \right\} \right).$$

8: Sketch of a sketch is a sketch

...or how many sketches would a sketch sketch, if a sketch could sketch sketches?

Lemma (Sketch property.)

$\mathbf{X} \subset \mathbf{Y}$: δ -sample.

$\mathbf{Y} \subset \mathbf{Z}$: δ' -sample.

$\implies \mathbf{X}$ is a $(\delta + \delta')$ -sample of \mathbf{Z} .

$\mathbf{Z} = \mathbf{P}$: Set of n points

$\exists \mathbf{Y} \subset \mathbf{Z}$: ϵ -Sample $|\mathbf{Y}| = \mathbf{O}\left(\frac{\log n}{\epsilon^2}\right)$.

$\exists \mathbf{X} \subset \mathbf{Y}$: ϵ -Sample $|\mathbf{X}| = \mathbf{O}\left(\frac{\log \log n}{\epsilon^2}\right)$.

\implies

\mathbf{X} is 2ϵ -sample of \mathbf{P} of size $\mathbf{O}\left(\frac{\log \log n}{\epsilon^2}\right)$.

Theorem [Vapnik and Chervonenkis, 1971]

Exists ϵ -sample of size $\mathbf{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)$.

8: Sketch of a sketch is a sketch

...or how many sketches would a sketch sketch, if a sketch could sketch sketches?

Lemma (Sketch property.)

$\mathbf{X} \subset \mathbf{Y}$: δ -sample.

$\mathbf{Y} \subset \mathbf{Z}$: δ' -sample.

$\implies \mathbf{X}$ is a $(\delta + \delta')$ -sample of \mathbf{Z} .

$\mathbf{Z} = \mathbf{P}$: Set of n points

$\exists \mathbf{Y} \subset \mathbf{Z}$: ε -Sample $|\mathbf{Y}| = \mathbf{O}\left(\frac{\log n}{\varepsilon^2}\right)$.

$\exists \mathbf{X} \subset \mathbf{Y}$: ε -Sample $|\mathbf{X}| = \mathbf{O}\left(\frac{\log \log n}{\varepsilon^2}\right)$.

\implies

\mathbf{X} is 2ε -sample of \mathbf{P} of size $\mathbf{O}\left(\frac{\log \log n}{\varepsilon^2}\right)$.

Theorem [Vapnik and Chervonenkis, 1971]

Exists ε -sample of size $\mathbf{O}\left(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon}\right)$.

8: Sketch of a sketch is a sketch

...or how many sketches would a sketch sketch, if a sketch could sketch sketches?

Lemma (Sketch property.)

$\mathbf{X} \subset \mathbf{Y}$: δ -sample.

$\mathbf{Y} \subset \mathbf{Z}$: δ' -sample.

$\implies \mathbf{X}$ is a $(\delta + \delta')$ -sample of \mathbf{Z} .

$\mathbf{Z} = \mathbf{P}$: Set of n points

$\exists \mathbf{Y} \subset \mathbf{Z}$: ϵ -Sample $|\mathbf{Y}| = \mathbf{O}\left(\frac{\log n}{\epsilon^2}\right)$.

$\exists \mathbf{X} \subset \mathbf{Y}$: ϵ -Sample $|\mathbf{X}| = \mathbf{O}\left(\frac{\log \log n}{\epsilon^2}\right)$.

\implies

\mathbf{X} is 2ϵ -sample of \mathbf{P} of size $\mathbf{O}\left(\frac{\log \log n}{\epsilon^2}\right)$.

Theorem [Vapnik and Chervonenkis, 1971]

Exists ϵ -sample of size $\mathbf{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)$.

8: Sketch of a sketch is a sketch

...or how many sketches would a sketch sketch, if a sketch could sketch sketches?

Lemma (Sketch property.)

$\mathbf{X} \subset \mathbf{Y}$: δ -sample.

$\mathbf{Y} \subset \mathbf{Z}$: δ' -sample.

$\implies \mathbf{X}$ is a $(\delta + \delta')$ -sample of \mathbf{Z} .

$\mathbf{Z} = \mathbf{P}$: Set of n points

$\exists \mathbf{Y} \subset \mathbf{Z}$: ϵ -Sample $|\mathbf{Y}| = \mathbf{O}\left(\frac{\log n}{\epsilon^2}\right)$.

$\exists \mathbf{X} \subset \mathbf{Y}$: ϵ -Sample $|\mathbf{X}| = \mathbf{O}\left(\frac{\log \log n}{\epsilon^2}\right)$.

\implies

\mathbf{X} is 2ϵ -sample of \mathbf{P} of size $\mathbf{O}\left(\frac{\log \log n}{\epsilon^2}\right)$.

Theorem [Vapnik and Chervonenkis, 1971]

Exists ϵ -sample of size $\mathbf{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)$.

9: Netting a lot with a little

or, getting less with even less

Smaller ϵ -samples are better.

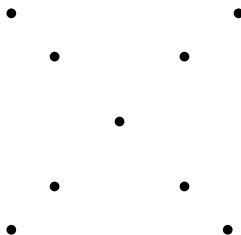
P: n points in the plane.

ϵ -Net

$\epsilon > 0$: Approx parameter.

S \subseteq **P** is ϵ -net if

$$\forall D \quad |\overline{P(D)}| \geq \epsilon \\ \implies S \cap D \neq \emptyset.$$



9: Netting a lot with a little

or, getting less with even less

Smaller ϵ -samples are better.

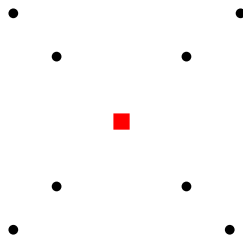
P: n points in the plane.

ϵ -Net

$\epsilon > 0$: Approx parameter.

S \subseteq **P** is ϵ -net if

$$\forall D \quad |\overline{P(D)}| \geq \epsilon \\ \implies S \cap D \neq \emptyset.$$



9: Netting a lot with a little

or, getting less with even less

Smaller ϵ -samples are better.

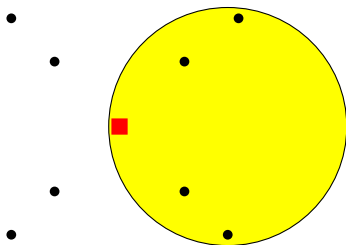
P: n points in the plane.

ϵ -Net

$\epsilon > 0$: Approx parameter.

$S \subseteq P$ is ϵ -net if

$$\forall D \quad |\overline{P(D)}| \geq \epsilon \\ \implies S \cap D \neq \emptyset.$$



10: ϵ -Sample vs. ϵ -net

P: n points in the plane.

$\epsilon > 0$: Approx parameter.

ϵ -Net

$S \subseteq P$ is ϵ -net if $\forall D \quad |\overline{P}(D)| \geq \epsilon \implies S \cap D \neq \emptyset.$

ϵ -Sample

$S \subseteq P$ is ϵ -sample if $\forall D \quad |\overline{P}(D) - \overline{S}(D)| < \epsilon.$

ϵ -sample is an ϵ -net.

Theorem (ϵ -net theorem)

[Haussler and Welzl, 1987]

Random sample of size $O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$ is an ϵ -net.

10: ϵ -Sample vs. ϵ -net

P: n points in the plane.

$\epsilon > 0$: Approx parameter.

ϵ -Net

$S \subseteq P$ is ϵ -net if $\forall D \quad |\overline{P(D)}| \geq \epsilon \implies S \cap D \neq \emptyset.$

ϵ -Sample

$S \subseteq P$ is ϵ -sample if $\forall D \quad |\overline{P(D)} - \overline{S(D)}| < \epsilon.$

ϵ -sample is an ϵ -net.

Theorem (ϵ -net theorem)

[Haussler and Welzl, 1987]

Random sample of size $O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$ is an ϵ -net.

10: ϵ -Sample vs. ϵ -net

P: n points in the plane.

$\epsilon > 0$: Approx parameter.

ϵ -Net

$S \subseteq P$ is ϵ -net if $\forall D \quad |\overline{P(D)}| \geq \epsilon \implies S \cap D \neq \emptyset.$

ϵ -Sample

$S \subseteq P$ is ϵ -sample if $\forall D \quad |\overline{P(D)} - \overline{S(D)}| < \epsilon.$

ϵ -sample is an ϵ -net.

Theorem (ϵ -net theorem)

[Haussler and Welzl, 1987]

Random sample of size $O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$ is an ϵ -net.

11: Hitting set via LP relaxation

ϵ -nets and Integrality gap.

LP relax' of hitting set.

$$\begin{aligned} \min \quad & \text{Opt} = \sum_{i=1}^n x_i \\ \text{s.t.} \quad & \sum_{p_i \in D_j} x_i \geq 1 \quad \forall j \\ & x_i \geq 0 \quad \forall i \end{aligned}$$

Lemma [Long, 2001]

In geometric settings there is an ϵ -net of size $O(K/\epsilon)$ iff the integrality gap is K .

ϵ -net of size...

$O(1/\epsilon \log 1/\epsilon) \implies O(\log \text{Opt})$ -approximation hitting-set

$O(1/\epsilon) \implies O(1)$ -approximation

12: On small ϵ -nets

Upper bounds

- $O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$: ϵ -net theorem.
- $O(1/\epsilon)$: halfplanes, halfspaces, pseudo-disks.
- $O\left(\frac{1}{\epsilon} \log \frac{U(1/\epsilon)}{1/\epsilon}\right)$

[Aronov, Ezra and Sharir, 2009]

$U(n)$: Union complexity of n shapes.

Lower bounds

- $\Omega\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$: Lower bound [Komlós et al., 1992].
- $\Omega\left(\frac{1}{\epsilon} w\left(\frac{1}{\epsilon}\right)\right)$: points and lines [Alon, 2010].

13: Other notions of samples

Name	Property $\forall \mathbf{D} \in \mathcal{F}$	Sample size
ϵ -net	$\bar{\mathbf{m}} = \bar{\mathbf{m}}(\mathbf{r}), \bar{\mathbf{s}} = \bar{\mathbf{s}}(\mathbf{r})$ $\bar{\mathbf{m}} \geq \epsilon \Rightarrow \bar{\mathbf{s}} > 0$	$O\left(\frac{\delta}{\epsilon} \log \frac{1}{\epsilon}\right)$
ϵ -sample	$ \bar{\mathbf{m}} - \bar{\mathbf{s}} \leq \epsilon$	$O\left(\frac{\delta}{\epsilon^2}\right)$
Sensitive ϵ -approx.	$ \bar{\mathbf{m}} - \bar{\mathbf{s}} \leq \frac{\epsilon}{2}(\sqrt{\bar{\mathbf{m}}} + \epsilon)$	$O\left(\frac{\delta}{\epsilon^2} \log \frac{1}{\epsilon}\right)$
Relative (ϵ, \mathbf{p}) -approx.	$\bar{\mathbf{m}} \geq \mathbf{p} \Rightarrow$ $(1 - \epsilon)\bar{\mathbf{m}} \leq \bar{\mathbf{s}} \leq (1 + \epsilon)\bar{\mathbf{m}}$	$O\left(\frac{\delta}{\epsilon^2 \mathbf{p}} \log \frac{1}{\mathbf{p}}\right)$

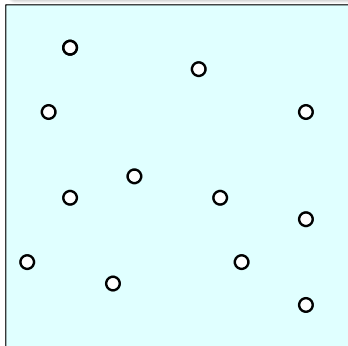
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E : # of edges of matching crossing h
Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

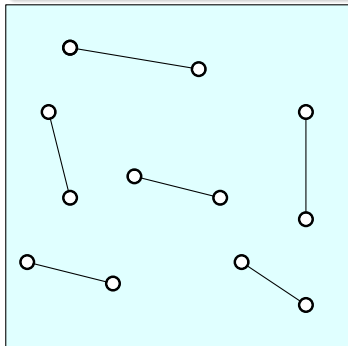
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E : # of edges of matching crossing h
Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

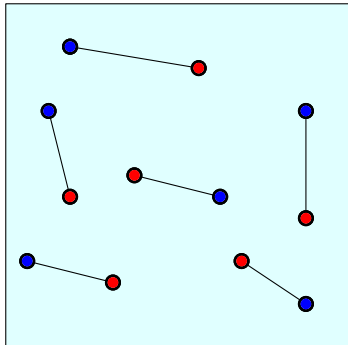
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E : # of edges of matching crossing h

Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

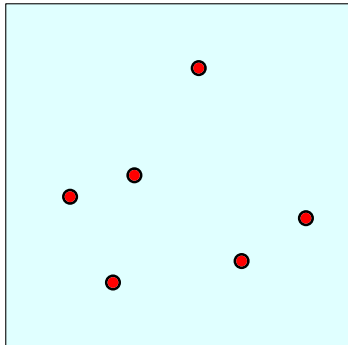
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E : # of edges of matching crossing h

Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

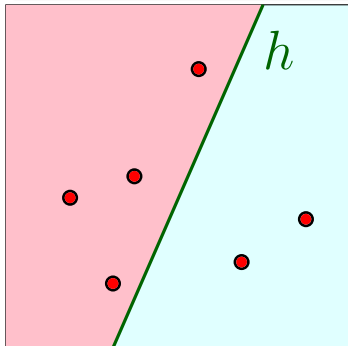
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E : # of edges of matching crossing h

Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

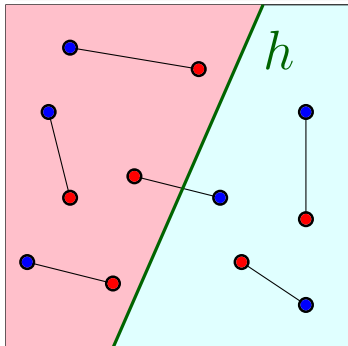
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to **P** using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E: # of edges of matching crossing **h**

Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

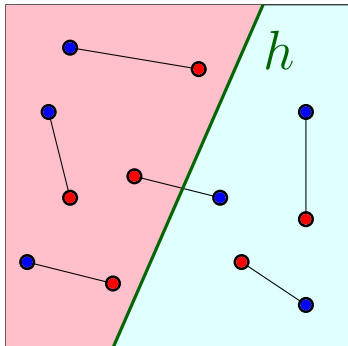
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E: # of edges of matching crossing h

Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

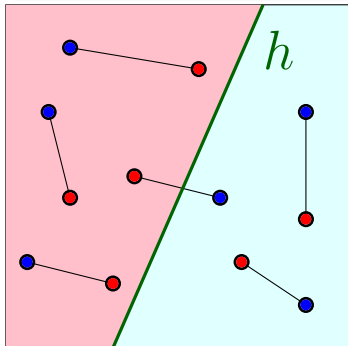
14: Getting a good sketch via discrepancy

Best sample using half the points?

Q: What is best approximation to \mathbf{P} using half the points?

Discrepancy

Color points by red/blue and take the red subset as approximation.



E: # of edges of matching crossing h

Result: Discrepancy $O(\sqrt{E \log n})$

Lemma

Red points form
 $\tilde{O}(1/\sqrt{n})$ -sample.

15: Spanning tree with low crossing number

[Welzl, 1992]

P: n points

T: spanning tree for **P**

Such that any line ℓ crosses

$O(\sqrt{n})$ edges of **T**

Proof: Uses reweighting.

15: Spanning tree with low crossing number

[Welzl, 1992]

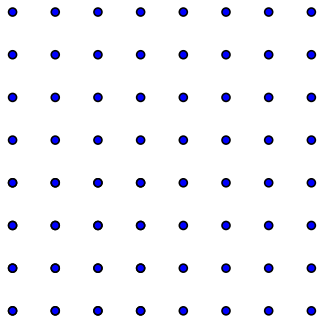
P: n points

T: spanning tree for **P**

Such that any line ℓ crosses

$O(\sqrt{n})$ edges of **T**

Proof: Uses reweighting.



15: Spanning tree with low crossing number

[Welzl, 1992]

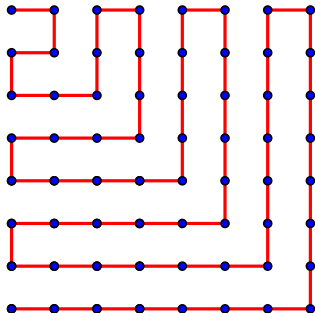
P: n points

T: spanning tree for **P**

Such that any line ℓ crosses

$O(\sqrt{n})$ edges of **T**

Proof: Uses reweighting.



15: Spanning tree with low crossing number

[Welzl, 1992]

P: n points

T: spanning tree for **P**

Such that any line ℓ crosses

$O(\sqrt{n})$ edges of **T**

Proof: Uses reweighting.

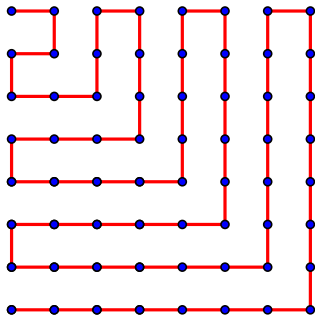
Any point set behaves like a grid.

\implies tour of **P** with

$cr(\mathbf{P}) = O(\sqrt{n})$.

\implies matching of **P** with

$cr(\mathbf{P}) = O(\sqrt{n})$.



16: Smaller ϵ -sample for halfplanes

Matching with low crossing number

$$n = |\mathbf{P}|$$

Plugging matching into disc. construction:

Theorem

\mathbf{P} : Set of n points in the plane.

One can compute coloring with discrepancy

$$O(\sqrt{\mathbf{E} \log n}) = \tilde{O}(n^{1/4}).$$

$\implies \tilde{O}(1/n^{3/4})$ -sample using $n/2$ points.

Theorem

For $(\text{points}, \text{halfplanes})$ there is ϵ -sample of size $\tilde{O}(1/\epsilon^{4/3})$.

17: Motivation: Another notion of a sketch

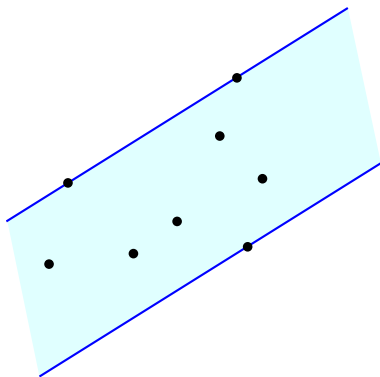
Coresets for Geometric Optimization

Compute width of a point set

Input: P - set points in \mathbb{R}^d .

Q: Find the width of P .

Find two parallel planes of min dist. that encloses P .



- Nasty problem for $d > 2$. (Exact algorithms are slow.)
- Approximation?

17: Motivation: Another notion of a sketch

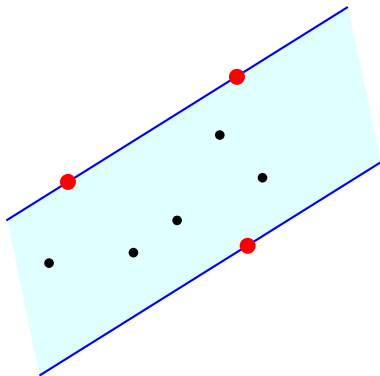
Coresets for Geometric Optimization

Compute width of a point set

Input: P - set points in \mathbb{R}^d .

Q: Find the width of P .

Find two parallel planes of min dist. that encloses P .



- Nasty problem for $d > 2$. (Exact algorithms are slow.)
- Approximation?

18: Motivation:

Geometric Optimization

- Find “quickly” a subset $Q \subseteq P$, s.t.,

$$\text{width}(Q) \geq (1 - \epsilon)\text{width}(P).$$

ϵ - approx parameter.

- Q - is a ϵ -coreset for P
(for the width function)
- Compute width of Q using naive algorithm.
- Running time $O(n + 1/\epsilon^{O(1)})$.

19: Coresets

- f : A monotone function defined over subsets of \mathbb{R}^d

$$S \subseteq T \Rightarrow f(S) \leq f(T).$$

- $P \subseteq \mathbb{R}^d$: input
- Q : Compute $f(P)$.
- $S \subseteq P$ is a ϵ -coreset for f if

$$f(S) \geq (1 - \epsilon)f(P).$$

20: Coresets - more general definition

- $f(\mathbf{P}, \mathbf{v})$: A monotone function defined over subsets of \mathbb{R}^d

$$\mathbf{S} \subseteq \mathbf{T} \Rightarrow f(\mathbf{S}, \mathbf{v}) \leq f(\mathbf{T}, \mathbf{v}).$$

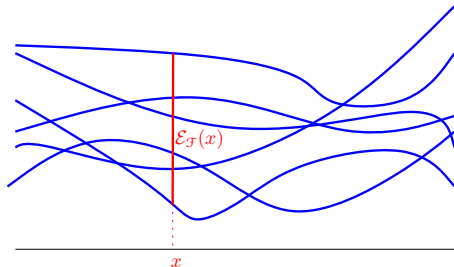
- $\mathbf{P} \subseteq \mathbb{R}^d$: input
- $\mathcal{S} \subseteq \mathbf{P}$ is a ϵ -coreset for \mathbf{f} if

$$\forall \mathbf{v} \quad f(\mathcal{S}, \mathbf{v}) \geq (1 - \epsilon)f(\mathbf{P}, \mathbf{v}).$$

21: Coresets - For family of functions

Vertical extent

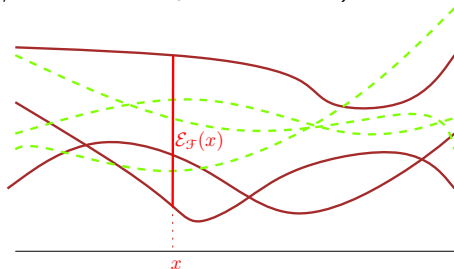
- \mathcal{F} - set of functions for \mathbb{R}^{d-1} to \mathbb{R} .
(i.e., surfaces in \mathbb{R}^d).
- $\mathcal{E}(\mathbf{x}) = \left(\max_{f \in \mathcal{F}} f(\mathbf{x}) \right) - \left(\min_{f \in \mathcal{F}} f(\mathbf{x}) \right)$
(Dist. upper/lower envelope of \mathcal{F} at \mathbf{x}).



21: Coresets - For family of functions

Vertical extent

- \mathcal{F} - set of functions for \mathbb{R}^{d-1} to \mathbb{R} .
(i.e., surfaces in \mathbb{R}^d).
- $\mathcal{E}(\mathbf{x}) = \left(\max_{f \in \mathcal{F}} f(\mathbf{x}) \right) - \left(\min_{f \in \mathcal{F}} f(\mathbf{x}) \right)$
(Dist. upper/lower envelope of \mathcal{F} at \mathbf{x}).



22: Coresets - For family of funcs

The result

[Agarwal et al., 2004]

\mathcal{F} : set of n functions

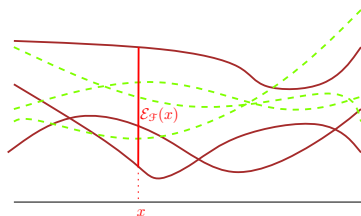
$\exists \mathcal{G} \subseteq \mathcal{F}$ s.t.

- \mathcal{G} is ϵ -coreset of \mathcal{F} :
 $\forall \mathbf{x} \quad \mathcal{E}_{\mathcal{G}}(\mathbf{x}) \geq (1 - \epsilon)\mathcal{E}_{\mathcal{F}}(\mathbf{x})$
- $|\mathcal{G}| = 1/\epsilon^{O(1)}$
- Construction time:
 $O(n + 1/\epsilon^{O(1)})$.

Functions in \mathcal{F} either:

- polynomials
- square root of polynomial.

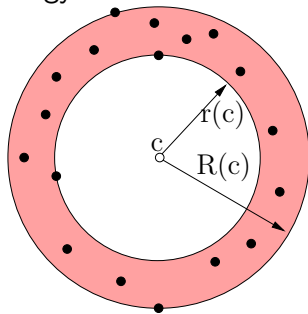
...c constant number of variables.



23: Application

Min-width annulus

- \mathbf{P} - set of n points in the plane.
- Find smallest annulus that covers \mathbf{P} .
- Motivation: meteorology.



Quality: $\mathbf{R(c) - r(c)}$

24: Min-width annulus

Reduction to extent

- $\mathbf{p}_i = (\mathbf{x}_i, \mathbf{y}_i) \in \mathbf{P}$ defines distance function:

$$f_i(\mathbf{c}) = f_i((\mathbf{x}, \mathbf{y})) = \|\mathbf{p}_i \mathbf{c}\| = \sqrt{(\mathbf{x}_i - \mathbf{x})^2 + (\mathbf{y}_i - \mathbf{y})^2}$$

$$\mathcal{F} = \{f_1, \dots, f_n\}$$

- Width of annulus at \mathbf{c} : $\mathcal{E}_{\mathcal{F}}(\mathbf{c}) = (\max_i f_i(\mathbf{c})) - (\min_j f_j(\mathbf{c}))$.
- Optimal annulus center: $\arg \min_{\mathbf{c} \in \mathbb{R}^2} \mathcal{E}_{\mathcal{F}}(\mathbf{c})$
- There is a coreset $\exists \mathcal{S} \subseteq \mathcal{F}$
 - (A) $\forall \mathbf{c} \quad \mathcal{E}_{\mathcal{S}}(\mathbf{c}) \geq (1 - \epsilon) \mathcal{E}_{\mathcal{F}}(\mathbf{c})$
 - (B) $|\mathcal{S}| = \mathbf{O}(1/\epsilon^2)$
- Compute shortest extent on coreset (brute force)
- Running time: $\mathbf{O}(n + 1/\epsilon^{\mathbf{O}(1)})$.

25: Other applications

- A generic/general algorithm.
- Other $(1 + \epsilon)$ -approximations:
 - Min radius cylinder.
 - Min width cylindrical shell.
 - Diameter.
 - Min enclosing ball.
 - Min volume bounding box.

All of the above for moving points.

- Handling outliers.
 - (A) $\mu(\mathbf{P}, \mathbf{k})$: min price of shape in \mathcal{F} cover all but \mathbf{k} points \mathbf{P} .
 - (B) $\exists \mathcal{S} \subseteq \mathbf{P}$ s.t. $\mu(\mathcal{S}, \mathbf{k}) \geq (1 - \epsilon)\mu(\mathbf{P}, \mathbf{k})$.
 $|\mathcal{S}| = \mathbf{O}(\mathbf{k}/\epsilon^d)$.
 - (C) Useful # of outliers small.

26: Sketches: Samples/Coresets/etc

Merging & sketching

Merge

- $\mathcal{S} \subseteq \mathbf{P}$ is ϵ -sketch
 - $\mathcal{T} \subseteq \mathbf{Q}$ is ϵ -sketch
- $\Rightarrow \mathcal{S} \cup \mathcal{T}$ is ϵ -sketch for $\mathbf{P} \cup \mathbf{Q}$.

Sketch

- $\mathbf{A} \subseteq \mathbf{B}$ is ϵ -sketch
 - $\mathbf{B} \subseteq \mathbf{C}$ is δ -sketch
- $\Rightarrow \mathbf{A}$ is $(\epsilon + \delta)$ -sketch for \mathbf{C} .

27: Sketches - streaming

- Points comes in a stream.
- Only small space is available.
- Maintain approx. to stream of points.

Result

If $\exists \epsilon$ -sketch for a problem of size $f(\epsilon)$, then one can solve problem on a stream using

$$O\left(f\left(\frac{\epsilon}{\log n}\right) \log n\right)$$

space.

n = number of elements in the stream.

28: Clustering

- **Clustering** = cover points with several shapes.
- Clustering price
 - L_∞ : Price of worst shape.
 - L_1 : Sum of prices.
- Example: **k** center problem
Cover point set with **k** balls of min max radius.
- Almost all clustering problems are NP-Complete.
- Extract sketch \Rightarrow cluster sketch.

29: Clustering - results

Fast algs. using coresets for all the following problems:

- **k**-center clustering (easy).
- **k**-center clustering for moving points.
- **k**-line center clustering
Cover points with **k** lines that min max dist of point to a line.
- **k**-median clustering
- **k**-means clustering.

30: Coresets in high dimensions

- \mathbf{P} - set of points in \mathbb{R}^d
 - $\mathcal{S} \subseteq \mathbf{P}$ is ε -coreset for min-ball if:
 - \mathbf{b} : ball in \mathbb{R}^d covering \mathcal{S}
 - $(1 + \varepsilon)\mathbf{b}$ covers \mathbf{P} .
 - Holds **for all** possible balls.
- ... Any coreset for min-ball has to be of size $\Omega(1/\varepsilon^{(d-1)/2})$.

31: Weak Sketches in high dimensions

- We still want to compute min enclosing ball.
- \mathbf{P} - set of points in \mathbb{R}^d
 - $\mathcal{S} \subseteq \mathbf{P}$ is ϵ -shellset for min-ball if:
 - \mathbf{b} : smallest ball in \mathbb{R}^d covering \mathcal{S}
 - $(1 + \epsilon)\mathbf{b}$ covers \mathbf{P} .

[Bădoiu and Clarkson, 2003]

\exists ϵ -shellset of size $O(1/\epsilon)$ of min-enclosing ball.
Computed in $O(dn/\epsilon + d/\epsilon^{O(1)})$ time.

32: Weak-sketches in high dimensions

Applications

“Fast” approximation to the following problems:

- Min enclosing ball with outliers
- Min enclosing cylinder (linear running time!)
- Min k -flat fitting
- Projective clustering with outliers.

The meta question...

To which optimization problems \exists short certificate of approximate optimal solution?






33: Conclusions

- Many problems in geometry possess small sketches.
- ... Far from having a unified theory.
- Open problems:

Approximate ϵ -sample

Given that \exists ϵ -sample of size t , compute a sample of size $O(t \log(1/\epsilon))$.

... How to compute/approximate most compact sketch?

-  Agarwal, P. K., Har-Peled, S., and Varadarajan, K. R. (2004).
Approximating extent measures of points.
J. Assoc. Comput. Mach., 51(4):606–635.
-  Alon, N. (2010).
A non-linear lower bound for planar epsilon-nets.
In *Proc. 51st Annu. IEEE Sympos. Found. Comput. Sci.*
-  Bădoiu, M. and Clarkson, K. L. (2003).
Optimal coresets for balls.
<http://cm.bell-labs.com/who/clarkson/coresets2.pdf>.
-  Haussler, D. and Welzl, E. (1987).
 ϵ -nets and simplex range queries.
Discrete Comput. Geom., 2:127–151.
-  Komlós, J., Pach, J., and Woeginger, G. (1992).
Almost tight bounds for ϵ -nets.
Discrete Comput. Geom., 7:163–173.



Li, Y., Long, P. M., and Srinivasan, A. (2001).
Improved bounds on the sample complexity of learning.
J. Comput. Syst. Sci., 62(3):516–527.



Long, P. M. (2001).
Using the pseudo-dimension to analyze approximation algorithms
for integer programming.
In Proc. 7th Workshop Algorithms Data Struct., volume 2125 of
Lecture Notes Comput. Sci., pages 26–37.



Vapnik, V. N. and Chervonenkis, A. Y. (1971).
On the uniform convergence of relative frequencies of events to
their probabilities.
Theory Probab. Appl., 16:264–280.



Welzl, E. (1992).
On spanning trees with low crossing numbers.

In *Data Structures and Efficient Algorithms, Final Report on the DFG Special Joint Initiative*, volume 594 of *Lect. Notes in Comp. Sci.*, pages 233–249. Springer-Verlag.