

Data approximation using low-crossing matchings

Mónika Csikós

joint work with Nabil Mustafa

INSTITUT DE RECHERCHE EN INFORMATIQUE FONDAMENTALE

How many people live within 100 km from a university?







How many people live within 100 km from a university?





How many people live within 100 km from a university?



Combinatorial data approximation







Combinatorial data approximation

Model with a set system: (X, \mathcal{S}) m subsets of X*n* elements

Goal: Pick half of the elements in X such that from each $S \in \mathcal{S}$,

'half' of the points in S are present in A.

Formally, $A \subset X$ of size $\frac{n}{2}$ that minimizes the discrepancy $\operatorname{disc}_{\mathcal{S}}(A) = \max_{S \in \mathcal{S}} \left| |S| - 2|A \cap S| \right|$



Random subset $A \subseteq X$ s.t. $x \in A$ with probabilit

set system $(X, \mathcal{S}), n = |X|, m = |\mathcal{S}|$

$$\operatorname{ty} \frac{1}{2} : \mathbb{E}|A \cap S| = \frac{|S|}{2}$$

Random subset $A \subseteq X$ s.t. $x \in A$ with probabilit

Chernoff-Hoeffding + union bound:

$$\mathbb{P}\left[\max_{S\in\mathcal{S}}|2|A\cap S|-|S|| > \delta\right] \le 2m \cdot \exp\left(-\frac{\delta^2}{2n}\right)$$

set system $(X, \mathcal{S}), n = |X|, m = |\mathcal{S}|$

$$\operatorname{ty} \frac{1}{2} : \mathbb{E}|A \cap S| = \frac{|S|}{2}$$

 $\implies \operatorname{disc}_{\mathscr{S}}(A) = O\left(\sqrt{n \ln m}\right)$

Recursively sample $A_1 \subseteq X$, $A_2 \subseteq A_1, A_3 \subseteq A_2, \ldots$, always halving the size

set system $(X, \mathcal{S}), n = |X|, m = |\mathcal{S}|$

Recursively sample $A_1 \subseteq X$, $A_2 \subseteq A_1, A_3 \subseteq A_2, \ldots$, always halving the size $|A_1| = \frac{n}{2}$ $2|A_1 \cap S| = |S| \pm O\left(\sqrt{n \ln m}\right)$ $|A_2| = \frac{n}{4}$ $4|A_2 \cap S| = |S| \pm O\left(\sqrt{2n\ln m}\right)$ • $|A_i| = \frac{n}{2^i} \qquad 2^i |A_i \cap S| = |S| \pm O\left(\sqrt{2^i n \ln m}\right)$

set system $(X, \mathcal{S}), n = |X|, m = |\mathcal{S}|$



Epsilon-approximations

 ε -approximation problem:

Given $\varepsilon \in (0,1)$, find smallest $A \subseteq X$ such that

Uniform sampling gives

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

$$\implies \text{A uniform sample of } O\left(\frac{\ln m}{\varepsilon^2}\right) \text{ points}$$

set system $(X, \mathcal{S}), n = |X|, m = |\mathcal{S}|$

$$\max_{S \in \mathcal{S}} \left| \frac{|S|}{|X|} - \frac{|A \cap S|}{|A|} \right| \le \varepsilon$$

pints is an ε -approximation with constant probability

Uniform sampling guarantees

set system (*X*, S), n = |X|, m = |S|

Arbitrary

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

ɛ-approximation size

ln m

 ε^2

set system (X, S)

X : set of *n* points in \mathbb{R}^2

each set in ${\mathcal S}$ is defined by a disk



set system (X, \mathcal{S})

X : set of *n* points in \mathbb{R}^2 each set in \mathcal{S} is defined by a disk $\implies |\mathcal{S}|_Y | < 2^{|Y|}$ even for |Y| = 4

 $\mathcal{S}|_{Y} = \{Y \cap S : S \in \mathcal{S}\}$



set system (X, \mathcal{S}) *X* : set of *n* points in \mathbb{R}^2 each set in \mathcal{S} is defined by a disk $\implies |\mathcal{S}|_{Y}| < 2^{|Y|}$ even for |Y| = 4

 $\mathcal{S}|_{Y} = \{Y \cap S : S \in \mathcal{S}\}$

but $\left| \mathcal{S} \right|_{Y} = 2^{3}$ possible for |Y| = 3



set system (X, \mathcal{S}) *X* : set of *n* points in \mathbb{R}^2 each set in \mathcal{S} is defined by a disk $\implies |\mathcal{S}|_{Y} | < 2^{|Y|}$ even for |Y| = 4 $\mathcal{S}|_{Y} = \{Y \cap S : S \in \mathcal{S}\}$

The VC-dimension of (X, \mathcal{S}) is the size of the large

- but $\left| \mathcal{S} \right|_{Y} = 2^{3}$ possible for |Y| = 3
- \implies VC-dimension of (X, \mathcal{S}) is 3.

est
$$Y \subseteq X$$
 s.t. $\left| \mathcal{S} \right|_Y = 2^{|Y|}$.





set system (X, \mathcal{S}) *X* : set of *n* points in \mathbb{R}^2 each set in \mathcal{S} is defined by a disk $\implies |\mathcal{S}|_{Y}| <$ but $\left| \mathcal{S} \right|_{Y} =$ $\mathcal{S}|_{Y} = \{Y \cap S : S \in \mathcal{S}\}$ \implies VC-dimens

The VC-dimension of (X, S) is the size of the large

$$d_{VC}(X, \mathcal{S}) = d \implies \left| \mathcal{S} \right|_{Y} = O\left(|Y|^{d} \right)$$

$$2^{|Y|}$$
 even for $|Y| = 4$
= 2^3 possible for $|Y| = 3$

sion of
$$(X, \mathcal{S})$$
 is 3.

est
$$Y \subseteq X$$
 s.t. $\left| \mathcal{S} \right|_Y = 2^{|Y|}$.

(Vapnik-Chervonenkis 1971, Sauer, Shelah 1972)



Uniform sampling guarantees

set system (*X*, S), n = |X|, m = |S|

Arbitrary

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

E-approximation size

 $\ln m$

 ε^2

VC-dim $\leq d$

 $dn \ln n$

 $\mathbf{1}$ \mathcal{E}^{\perp}

(Li-Long-Srinivasan 2001)

X : set of *n* points in \mathbb{R}^2 each set in *S* is defined by a disk

X : set of *n* points in \mathbb{R}^2 each set in *S* is defined by a disk

Theorem (Welzl 1988, Chazelle-Welzl 1989)

X has a perfect matching such that any disk c



crosses
$$O\left(\sqrt{n}\right)$$
 edges.

X : set of *n* points in \mathbb{R}^2 each set in ${\mathcal S}$ is defined by a disk

Theorem (Welzl 1988, Chazelle-Welzl 1989)

X has a perfect matching such that any disk c

Idea: Take one end-point of each matching edge randomly \implies only $O\left(\sqrt{n}\right)$ random variables for each disk $\mathbb{P}\left[\max_{S\in\mathcal{S}}\left|2|A\cap S|-|S|\right| > \delta\right] \le m$

 \implies disc_S(A) = O($n^{1/4}\sqrt{\ln n}$) with constant probability

crosses
$$O\left(\sqrt{n}\right)$$
 edges.

$$\cdot \exp\left(-\frac{\delta^2}{2\sqrt{n}}\right)$$

- (Matoušek-Welzl-Wernish 1991)



X : set of *n* points in \mathbb{R}^2 each set in ${\mathcal S}$ is defined by a disk

Theorem (Welzl 1988, Chazelle-Welzl 1989)

X has a perfect matching such that any disk c

Idea: Take one end-point of each matching edge randomly \implies only $O\left(\sqrt{n}\right)$ random variables for each disk $\mathbb{P}\left[\max_{S\in\mathcal{S}}|2|A\cap S|-|S|| > \delta\right] \le m$

 \implies disc_S(A) = $O\left(n^{1/4}\sqrt{\ln n}\right)$ with constant probability

crosses
$$O\left(\sqrt{n}\right)$$
 edges.

$$\cdot \exp\left(-\frac{\delta^2}{2\sqrt{n}}\right)$$

- (Matoušek-Welzl-Wernish 1991)

X : *n* points in \mathbb{R}^d \mathcal{S} : balls in \mathbb{R}^d

 $O(n^{1-1/d})$

$$O\left(\sqrt{dn^{1-1/d}\ln n}\right)$$

Approximation bounds

X : *n* points in \mathbb{R}^d , \mathcal{S} : subsets induced by balls

Uniform sampling

Arbitrary

Discrepancy error

 $n \ln m$

 ε -approximation size

 $\ln m$

 ε^2

Structural properties + non-uniform sampling

VC-dim $\leq d$

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

 $dn^{1-1/d}\ln n$

$$\left(\frac{d\ln\frac{1}{\varepsilon}}{\varepsilon^2}\right)^{1-1/d}$$

 ε^2



Approximation bounds

X : *n* points in \mathbb{R}^d , \mathcal{S} : subsets induced by balls

Uniform sampling

Arbitrary

Discrepancy error

 $n \ln m$

 ε -approximation size

 $\ln m$

 ε^2

Structural properties + non-uniform sampling

VC-dim $\leq d$

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

 ε^2

$$\sqrt{dn^{1-1/d}\ln n}$$

$$\left(\frac{d\ln\frac{1}{\varepsilon}}{\varepsilon^2}\right)^{1-1/d}$$

Need a matching to sample





Iterative reweighing method

Maintain weights on balls

- 1. initially, each ball has weight = 1
- 2. for i = 1, ..., n/2
 - select an edge e_i crossing balls of minimum total weight
 - double the weight of each balls crossing e_i
 - remove all edges incident to e_i

Return $e_1, ..., e_{n/2}$







Iterative reweighing method

Maintain weights on balls

- 1. initially, each ball has weight = 1
- 2. for i = 1, ..., n/2
 - select an edge e_i crossing balls of minimum total weight
 - double the weight of each balls crossing e_i
 - remove all edges incident to e_i

Return $e_1, ..., e_{n/2}$



 $O(n^2m)$ time





Iterative reweighing method

Maintain weights on balls

- 1. initially, each ball has weight = 1
- 2. for i = 1, ..., n/2
 - select an edge e_i crossing balls of minimum total weight
 - double the weight of each balls crossing e_i
 - remove all edges incident to e_i

Return $e_1, ..., e_{n/2}$





(Chazelle-Welzl 1989)



Maintain weights on edges and balls

1. initially, each ball and edge has weight = 1

2.

- sample an edge e_i according to the current edge-weights
- sample a ball S_i according to the current range-weights
- double the weight of each ball crossing e_i
- halve the weight of each edge crossing S_i

edge-weights nge-weights





Maintain weights on edges and balls

1. initially, each ball and edge has weight = 1

2.

- sample an edge e_i according to the current edge-weights
- sample a ball S_i according to the current range-weights
- double the weight of each ball crossing e_i
- halve the weight of each edge crossing S_i

edge-weights nge-weights



Maintain weights on edges and balls

1. initially, each ball and edge has weight = 1

2.

- sample an edge e_i according to the current edge-weights
- sample a ball S_i according to the current range-weights
- double the weight of each ball crossing e_i
- halve the weight of each edge crossing S_i

edge-weights nge-weights





Maintain weights on edges and balls

1. initially, each ball and edge has weight = 1

2.

- sample an edge e_i according to the current edge-weights
- sample a ball S_i according to the current range-weights
- double the weight of each ball crossing e_i
- halve the weight of each edge crossing S_i
- add e_i to the matching and remove all edges incident to e_i

Maintain weights on edges and balls

- 1. initially, each ball and edge has weight = 1
- 2. for i = 1, ..., n/4
 - sample an edge e_i according to the current edge-weights
 - sample a ball S_i according to the current range-weights
 - double the weight of each ball crossing e_i
 - halve the weight of each edge crossing S_i
 - add e_i to the matching and remove all edges incident to e_i

3. recurse on uncovered points

E crossing number of $e_1, \dots, e_{n/4} \lesssim n^{1-1/d}$

 $+ \ln m$



Maintain weights on edges and balls

- 1. initially, each ball and edge has weight = 1
- 2. for i = 1, ..., n/4
 - sample an edge e_i according to the current edge-weights
 - sample a ball S_i according to the current range-weights
 - double the weight of each ball crossing e_i
 - halve the weight of each edge crossing S_i
 - add e_i to the matching and remove all edges incident to e_i

3. recurse on uncovered points

$$O(m+n^2)$$
 time

E crossing number of $e_1, \dots, e_{n/4} \lesssim n^{1-1/d}$

 $+ \ln m$



Maintain weights on edges and balls

- 1. initially, each ball and edge has weight = 1
- 2. for i = 1, ..., n/4
 - sample an edge e_i according to the current edge-weights
 - sample a ball S_i according to the current range-weights
 - double the weight of each ball crossing e_i
 - halve the weight of each edge crossing S_i
 - add e_i to the matching and remove all edges incident to e_i

3. recurse on uncovered points





Our method

Maintain weights on edges and balls

0. sample $E_1, \ldots, E_{n/4} \subset E$ and $\mathcal{S}_1, \ldots, \mathcal{S}_{n/4} \subset \mathcal{S}$

1. initially, each range and edge has weight = 1

2. for i = 1, ..., n/4

- sample an edge e_i according to the current edge-weights
- sample a ball S_i according to the current range-weights
- double the weight of each ball in S_i crossing e_i
- halve the weight of each edge in E_i crossing S_i
- add e_i to the matching and remove all edges incident to e_i

3. recurse on uncovered points

with
$$\mathbb{E}\left[|\mathcal{S}_i|\right] = \frac{m \ln m}{n^{1-1/d}}$$
 and $\mathbb{E}\left[|E_i|\right] = n^{1+1/d} \ln m$

Matching with crossing number $O(n^{1-1/d})$ in time $\tilde{O}(mn^{1/d} + n^{2+1/d})$



Matchings – \mathbb{R}^2

5000 points on 10 co-centric circles

Our method, Ranges: half-planes











Matchings – \mathbb{R}^2

5000 points picked uniformly at random from 5000 grid cells

Our method, Ranges: half-planes







Our method, Ranges: disks

Random Matching



Matchings – \mathbb{R}^2

5000 points picked uniformly at random from 5000 grid cells





Approximations — \mathbb{R}^2

Input: 10000 points

Our method



Input: 10000 points

Random sampling



