

Sampling in Geometric and Combinatorial Set Systems

NABIL H. MUSTAFA



ESIEE
PARIS

UNIVERSITÉ —
— **PARIS-EST**

OUTLINE

Configurations

graphs, sets, k -tuples, random, ...



Approximations

samples, coresets, kernels, sketches, ...



Applications

optimisation, graph algorithms, learning, ...

APPROXIMATIONS

n elements \swarrow m sets $= O(n^d)$
 set system (X, \mathcal{R})

$A \subseteq X$

$S \in \mathcal{R}$

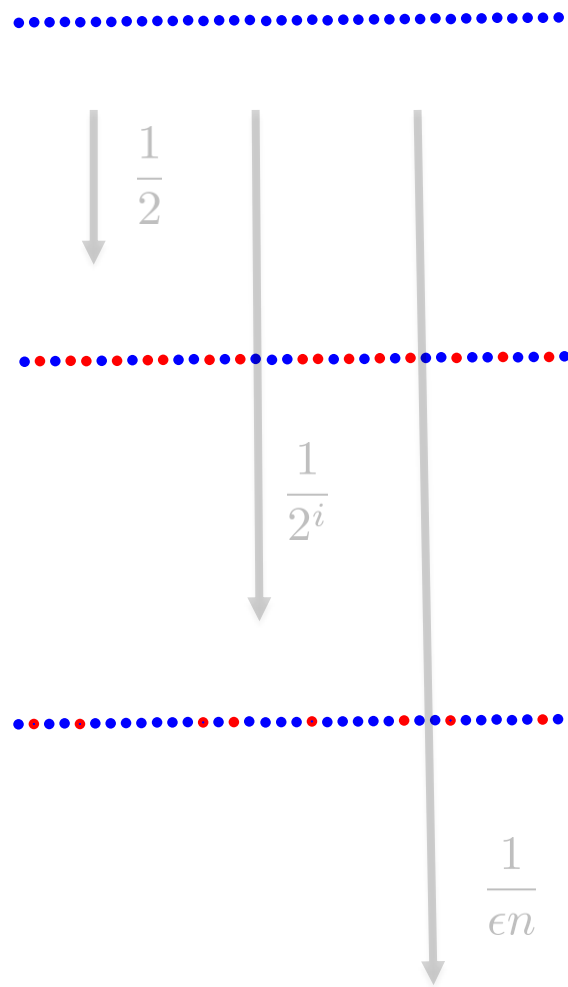
$ A $	$\mathbf{E}[A \cap S]$	discrepancy
$t = \frac{n}{2}$	$\frac{ S }{2}$	$\frac{ S }{2} \pm \boxed{\text{error}}$

i steps $\frac{|S|}{2^i} \pm \epsilon t$

$\boxed{t} = \frac{n}{2^i}$	$\frac{ S }{2^i}$	ϵ-approximations
		$\frac{ S t}{n} \pm \epsilon t$

$\log(\epsilon n)$ steps, $\epsilon > 0$ a given parameter

$t = \frac{1}{\epsilon}$	$\frac{ S }{\epsilon n}$	
--------------------------	--------------------------	--



APPROXIMATIONS

n elements $\quad m$ sets $= O(n^d)$
 set system (X, \mathcal{R})

$A \subseteq X$

$S \in \mathcal{R}$

$|A|$

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

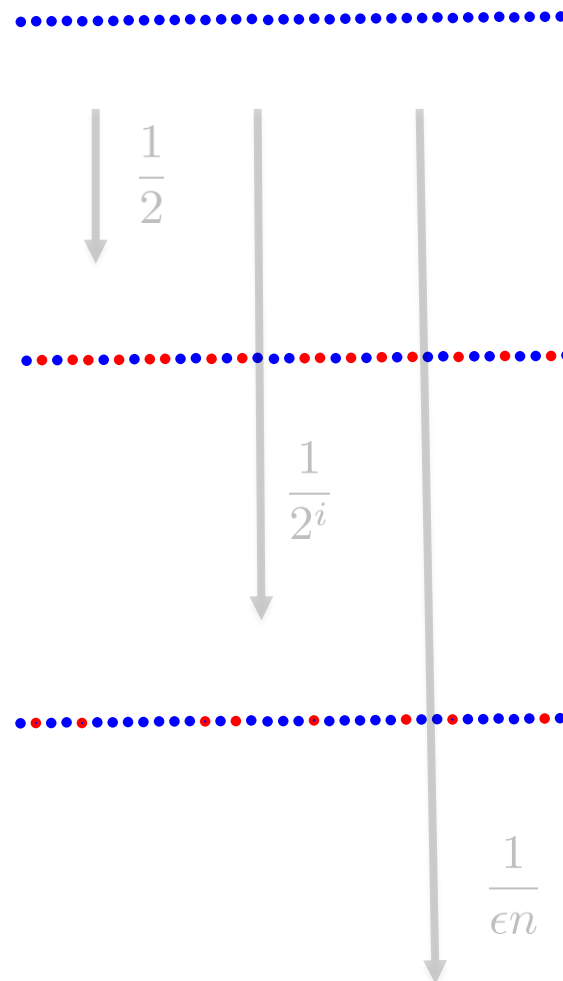
$t = \frac{n}{2}$	$\frac{ S }{2}$	$\frac{ S }{2} \pm \boxed{\text{error}}$ discrepancy
-------------------	-----------------	--

i steps

$\boxed{t} = \frac{n}{2^i}$	$\frac{ S }{2^i}$	$\frac{ S t}{n} \pm \epsilon t$ ϵ-approximations
-----------------------------	-------------------	---

$\log(\epsilon n)$ steps, $\epsilon > 0$ a given parameter

$\boxed{t} = \frac{1}{\epsilon}$	$\frac{ S }{\epsilon n}$	≥ 1 if $ S \geq \epsilon n$ ϵ-nets
----------------------------------	--------------------------	--



RANDOM SAMPLING FOR APPROXIMATIONS

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

A : a uniform random sample of X of size t .

$$\mathbf{E}[|A \cap S|] = |S| \cdot \frac{t}{n}$$

Chernoff's bound: For any $\eta > 0$ and $S \in \mathcal{R}$

$$\Pr \left[|A \cap S| \notin \frac{|S|t}{n} \pm \eta \right] \leq 2 \exp \left(-\frac{\eta^2 n}{2|S|t + \eta n} \right) \leq 2 \exp \left(-\frac{\eta^2}{3t} \right).$$

discrepancy

$$t = \frac{n}{2} \quad \frac{|S|}{2} \pm \text{error}$$

$$t = \frac{n}{2}$$

probability of discrepancy
being at least η

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

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

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

$$\mathcal{O} \left(\sqrt{n \ln m} \right)$$

RANDOM SAMPLING FOR APPROXIMATIONS

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

A : a uniform random sample of X of size t .

$$\mathbf{E}[|A \cap S|] = |S| \cdot \frac{t}{n}$$

Chernoff's bound: For any $\eta > 0$ and $S \in \mathcal{R}$

$$\Pr \left[|A \cap S| \notin \frac{|S|t}{n} \pm \eta \right] \leq 2 \exp \left(-\frac{\eta^2 n}{2|S|t + \eta n} \right) \leq 2 \exp \left(-\frac{\eta^2}{3t} \right).$$

discrepancy

$$t = \frac{n}{2} \quad \frac{|S|}{2} \pm \text{error}$$

$$t = \frac{n}{2}$$

probability of discrepancy
being at least η

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

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

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

$$O \left(\sqrt{n \ln m} \right)$$

ϵ -approximations

$$\frac{|S|t}{n} \pm \epsilon t$$

$$\eta = \epsilon t \quad t = \Theta \left(\frac{1}{\epsilon^2} \ln \frac{m}{\gamma} \right)$$

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

$$\leq m \cdot \exp \left(-\Theta \left(\epsilon^2 t \right) \right)$$

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

$$O \left(\frac{1}{\epsilon^2} \ln m \right)$$

RANDOM SAMPLING FOR APPROXIMATIONS

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

A : a uniform random sample of X of size t .

$$\mathbf{E}[|A \cap S|] = |S| \cdot \frac{t}{n}$$

Chernoff's bound: For any $\eta > 0$ and $S \in \mathcal{R}$

$$\Pr \left[|A \cap S| \notin \frac{|S|t}{n} \pm \eta \right] \leq 2 \exp \left(-\frac{\eta^2 n}{2|S|t + \eta n} \right) \leq 2 \exp \left(-\frac{\eta^2}{3t} \right).$$

discrepancy

$$t = \frac{n}{2} \quad \frac{|S|}{2} \pm \text{error}$$

$$t = \frac{n}{2}$$

probability of discrepancy being at least η

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

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

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

$$O \left(\sqrt{n \ln m} \right)$$

ϵ -approximations

$$\frac{|S|t}{n} \pm \epsilon t$$

$$\eta = \epsilon t \quad t = \Theta \left(\frac{1}{\epsilon^2} \ln \frac{m}{\gamma} \right)$$

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

$$\leq m \cdot \exp \left(-\Theta \left(\epsilon^2 t \right) \right)$$

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

$$O \left(\frac{1}{\epsilon^2} \ln m \right)$$

ϵ -nets

$$\geq 1 \quad \text{for } |S| \geq \epsilon n$$

$$t = \Theta \left(\frac{1}{\epsilon} \ln \frac{m}{\gamma} \right)$$

probability of not hitting some set $S \in \mathcal{R}$

$$\leq m \cdot (1 - \epsilon)^t \leq m \cdot \exp(-\epsilon t)$$

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

$$O \left(\frac{1}{\epsilon} \ln m \right)$$

RANDOM SAMPLING FOR APPROXIMATIONS

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

A : a uniform random sample of X of size t .

$$\mathbf{E}[|A \cap S|] = |S| \cdot \frac{t}{n}$$

Random Sampling

Chernoff's bound: For any $\eta > 0$ and $S \in \mathcal{R}$

$$\Pr \left[|A \cap S| \notin \frac{|S|t}{n} \pm \eta \right] \leq 2 \exp \left(-\frac{\eta^2 n}{2|S|t + \eta n} \right) \leq 2 \exp \left(-\frac{\eta^2}{3t} \right).$$

discrepancy

$$t = \frac{n}{2} \quad \frac{|S|}{2} \pm \text{error}$$

$$t = \frac{n}{2}$$

probability of discrepancy
being at least η

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

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

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

$$O \left(\sqrt{n \ln m} \right)$$

ϵ -approximations

$$\frac{|S|t}{n} \pm \epsilon t$$

$$\eta = \epsilon t \quad t = \Theta \left(\frac{1}{\epsilon^2} \ln \frac{m}{\gamma} \right)$$

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

$$\leq m \cdot \exp \left(-\Theta \left(\epsilon^2 t \right) \right)$$

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

$$O \left(\frac{1}{\epsilon^2} \ln m \right)$$

ϵ -nets

$$\geq 1 \quad \text{for } |S| \geq \epsilon n$$

$$t = \Theta \left(\frac{1}{\epsilon} \ln \frac{m}{\gamma} \right)$$

probability of not hitting
some set $S \in \mathcal{R}$

$$\leq m \cdot (1 - \epsilon)^t \leq m \cdot \exp(-\epsilon t)$$

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

$$O \left(\frac{1}{\epsilon} \ln m \right)$$

APPLICATION

ϵ -nets

OPTIMISATION

Minimum hitting set problem on input (P, \mathcal{R})

[Vazirani 2003]

Algorithm: **GREEDY**

$N = \emptyset$

While N not a hitting set for \mathcal{R}

 Add $q \in P$ that hits maximum new sets, to N

Return N

Algorithm: **LINEAR PROGRAM**

Solve LP

N : random sample of P w.r.t weights $x_p \cdot \log m$

Return N

$O(\text{OPT} \cdot \log m)$

Minimize $\sum_{p \in P} x_p \leq \text{OPT}$

subject to

(C1) $\sum_{p \in R} x_p \geq 1 \quad \forall R \in \mathcal{R},$

(C2) $0 \leq x_p \leq 1 \quad \forall p \in P.$

$\epsilon = \frac{1}{\text{OPT}} \rightarrow O(\text{OPT} \cdot \log m)$

Claim : If (P, \mathcal{R}) has ϵ -nets of size $\frac{1}{\epsilon} \cdot f\left(\frac{1}{\epsilon}\right)$
then $f(\text{OPT})$ -approximation

[Long 2001]

OPTIMISATION

Minimum hitting set problem on input (P, \mathcal{R})

Algorithm: RANDOMIZED GREEDY + LP

Solve LP

N : random sample of P w.r.t weights $x_p \cdot \log m$ ~~$x_p \cdot \log m$~~ $x_p \cdot S$

While N not a hitting set for \mathcal{R} $S = 1, 2, \dots, \log m$

R : any set of \mathcal{R} not hit by N

q : an element randomly sampled from R according to weights $\{x_p : p \in R\}$

$N = N \cup \{q\}$

Return N

[M. 2019]

Claim: optimal approximation bounds (within constant factors)
for most well-studied geometric systems

runs for expected OPT iterations

Minimize $\sum_{p \in P} x_p$

subject to

(C1) $\sum_{p \in R} x_p \geq 1 \quad \forall R \in \mathcal{R},$

(C2) $0 \leq x_p \leq 1 \quad \forall p \in P.$

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 with constant probability.

Surprising Theorem: A uniform random sample of X of size $\Theta\left(\frac{1}{\epsilon^2} \ln \frac{1}{\epsilon}\right)$ is an ϵ -approximation with constant probability, for a **locally nice system**

[Vapnik, Chervonenkis 71]

‘locally nice’ set system

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

number of subsets on $Y \subseteq X: O(|Y|^d)$

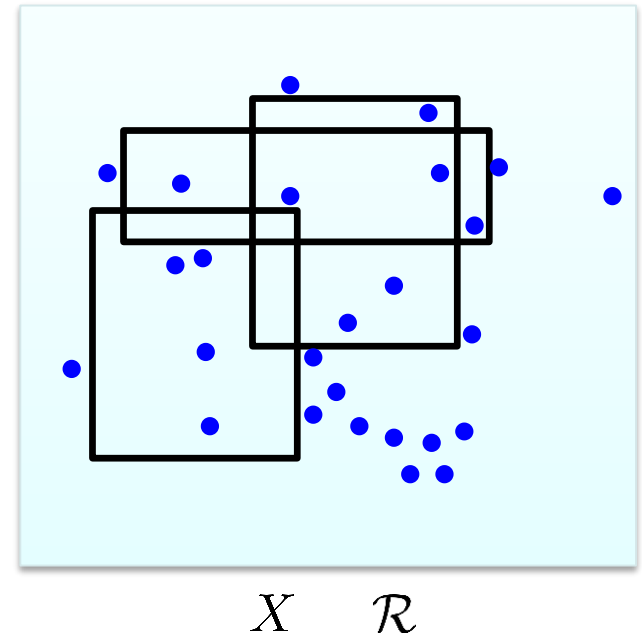
combinatorially

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

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

a constant d such that

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



SAMPLING IN NICE SYSTEMS

ϵ -approximations: A uniform random sample of X of size $\Theta\left(\frac{d}{\epsilon^2} \ln \frac{1}{\epsilon}\right)$ is an ϵ -approximation of \mathcal{R} , with constant probability, for nice systems. [Vapnik, Chervonenkis 71]

typically computational learning theory uses more complicated technique called *symmetrization*

Previous bound : a uniform random sample of size $t = \Theta\left(\frac{1}{\epsilon^2} \ln m\right)$ is an ϵ -approximation of \mathcal{R} .

$$X \quad |X| = n$$

$\frac{\epsilon}{2}$ -approximation
of \mathcal{R}



$$|A'| = O\left(\frac{1}{\epsilon^2} \ln n^d\right)$$

Random Sampling

A'

$\frac{\epsilon}{2}$ -approximation
of $\mathcal{R}|_{A'}$



$$|A| = O\left(\frac{1}{\epsilon^2} \ln |A'|^d\right) = O\left(\frac{d}{\epsilon^2} \ln \left(\frac{1}{\epsilon^2} \ln n^d\right)\right) = O\left(\frac{d}{\epsilon^2} \ln \frac{1}{\epsilon^2} + \frac{d}{\epsilon^2} \ln \ln n^d\right)$$

A

ϵ -approximation of \mathcal{R}

$$T(\epsilon) = O\left(\frac{1}{\epsilon^2} \ln T\left(\frac{\epsilon}{2}\right)^d\right)$$



$$\Theta\left(\frac{d}{\epsilon^2} \ln \frac{1}{\epsilon}\right)$$

[Csikos, M., 2020]

ϵ -nets: A uniform random sample of X of size $O\left(\frac{d}{\epsilon} \log \frac{1}{\epsilon}\right)$ is an ϵ -net with constant probability

[Haussler, Welzl 87]

VC DIMENSION

VC Dimension: A classical measure of complexity of set systems (X, \mathcal{R})

size of the largest set for which all subsets are possible; that is,

$$|\mathcal{R}|_Y = 2^{|Y|}$$

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

Sauer-Shelah lemma

related to the dimension of the Euclidean space

$d + 1$ for half-spaces in \mathbb{R}^d

$d + 1$ for balls in \mathbb{R}^d

$$|\mathcal{R}|_Y = O(|Y|^d)$$

an important case: VC dimension of union of k balls, half-spaces in \mathbb{R}^d

key parameter behind certain high-dimensional clustering algorithms

$O(d \cdot k \log k)$ [Blumer, Ehrenfeucht, Haussler, Warmuth '89]

$\Omega(d \cdot k \log k)$ { probabilistic construction: take a random set system [Eisenstat, Angluin, 2007]
a different construction for half-spaces, balls

[Csikos, Kupavskii, M., 2019]

alterations: pay too much if we sample large-enough to eliminate *all* 'bad events'

probability, statistics, learning

the biggest risk is
to not take any risk!

n elements
 m subsets
 d dimension

Ideal

Random
Arbitrary

VC dimension

Sampling

Sampling
+
Combinatorics

Discrepancy

0

$\sqrt{n \ln m}$

$\sqrt{dn \ln n}$



Approximations

$\frac{1}{\epsilon}$

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

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



Nets

$\frac{1}{\epsilon}$

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

$\frac{d}{\epsilon} \ln \frac{1}{\epsilon}$

chaining

random set system

[Komlos, Pach, Woeginger 1992]

[Talagrand, 1994]

[Li, Long, Srinivasan, 2001]

EPSILON-NETS FOR ABSTRACT SET SYSTEMS

ϵ -nets: A uniform random sample of X of size $O\left(\frac{d}{\epsilon} \log \frac{1}{\epsilon}\right)$ is an ϵ -net with constant probability

NON-OPTIMAL

Throughout the 1990s and the 2000s

several new techniques developed for $o\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$ sized nets

[Clarkson, Varadarajan 2006] [Pyrga, Ray 2008] [Varadarajan 2008]

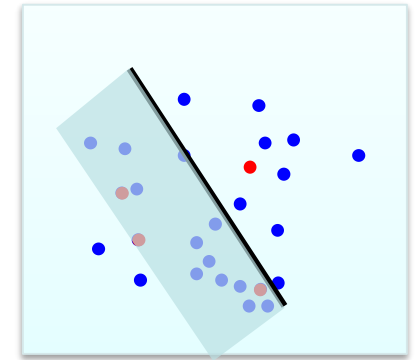
[Varadarajan 2009] [Aronov, Ezra, Sharir 2010] [Chan et al. 2012]

[Har-Peled, Kaplan, Sharir, Smorodinsky 2014] [Dutta, Ghosh, M. 2018]

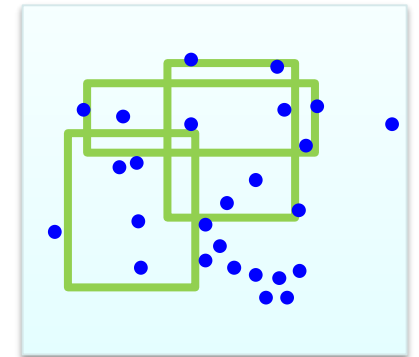
VC-dimension not fine enough

shortest paths in planar graphs

$O\left(\frac{1}{\epsilon}\right)$ -sized nets



$O\left(\frac{1}{\epsilon}\right)$ -sized nets



$O\left(\frac{1}{\epsilon} \log \log \frac{1}{\epsilon}\right)$ -sized nets

THE STRENGTHENED ϵ -NET THEOREM

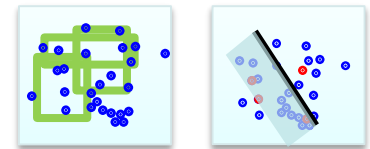
INSIGHT



(X, \mathcal{R}) has *shallow-cell complexity* $\varphi(\cdot, \cdot)$ if for any $Y \subseteq X$, and integer k
number of sets in $\mathcal{R}|_Y$ of size at most k is $O(|Y| \cdot \varphi(|Y|, k))$

Strengthened Epsilon-Net theorem: Let (X, \mathcal{R}) be a set system with shallow-cell complexity $\varphi(\cdot, \cdot)$. Then there exists an ϵ -net of size

$$O\left(\frac{d}{\epsilon} + \frac{1}{\epsilon} \log \varphi_{\mathcal{F}}\left(\frac{8d}{\epsilon}, 24d\right)\right)$$



[Varadarajan 2009]

[Chan et al. 2012]

[Dutta, Ghosh, M. 2018]

algorithm requires computing additional structures

optimal: probabilistic construction
 [Kupavskii, M., Pach 2017]

AN ALGORITHM

General Net-Finder $((\mathbf{X}, \mathcal{F}), \epsilon > 0)$.

N : pick a uniform random subset of X

while *there exists a set $S \in \mathcal{F}$ not hit by N* **do**

└ Add $\Theta(1)$ uniformly chosen random elements of S to N .

return N .

chaining + alterations

gives all known bounds

AN ALGORITHM

General Net-Finder $((\mathbf{X}, \mathcal{F}), \epsilon > 0)$.

N : pick a uniform random subset of X

while there exists a set $S \in \mathcal{F}$ not hit by N **do**

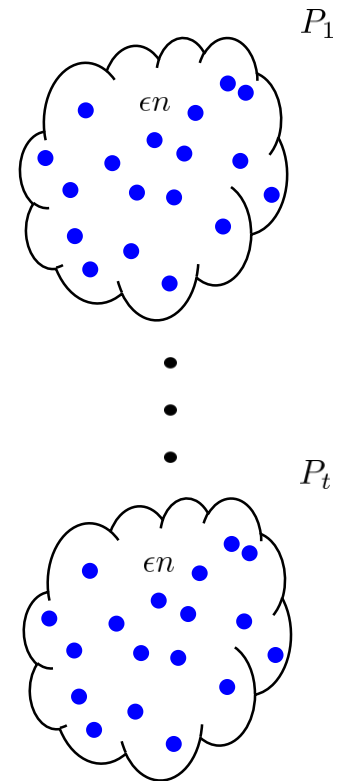
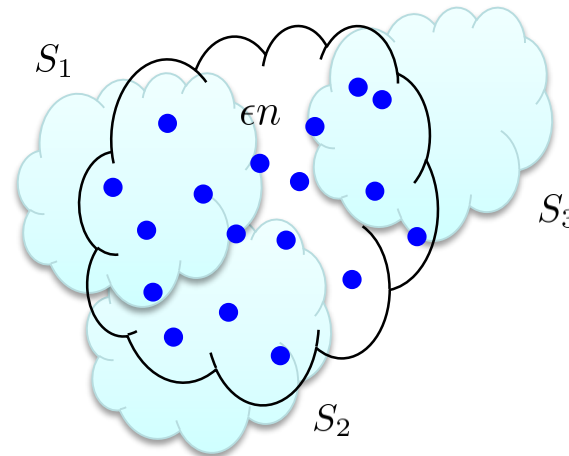
\perp Add $\Theta(1)$ uniformly chosen random elements of S to N .

return N .

[M. 2019]

1. **clustering.** there are a few key sets of \mathcal{F} $O\left(\frac{1}{\epsilon} \varphi_{\mathcal{F}}\left(\frac{d}{\epsilon}, d\right)\right)$
2. **reduction.** the initial random sample reduces it even further
3. **potential function.** each iteration makes progress on one of them

probabilistic charging argument



[Haussler 1995]

n elements
 m subsets
 d dimension

Ideal

Random

VC dimension

Sampling

Sampling
 +
 Combinatorics

Discrepancy

0

$\sqrt{n \ln m}$

$\sqrt{n \ln m}$

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



Approximations

$\frac{1}{\epsilon}$

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

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

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



Nets

$\frac{1}{\epsilon}$

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

$\frac{d}{\epsilon} \ln \frac{1}{\epsilon}$

$O\left(\frac{1}{\epsilon} \log \varphi_{\mathcal{F}}\left(\frac{16d}{\epsilon}, 64d\right)\right)$

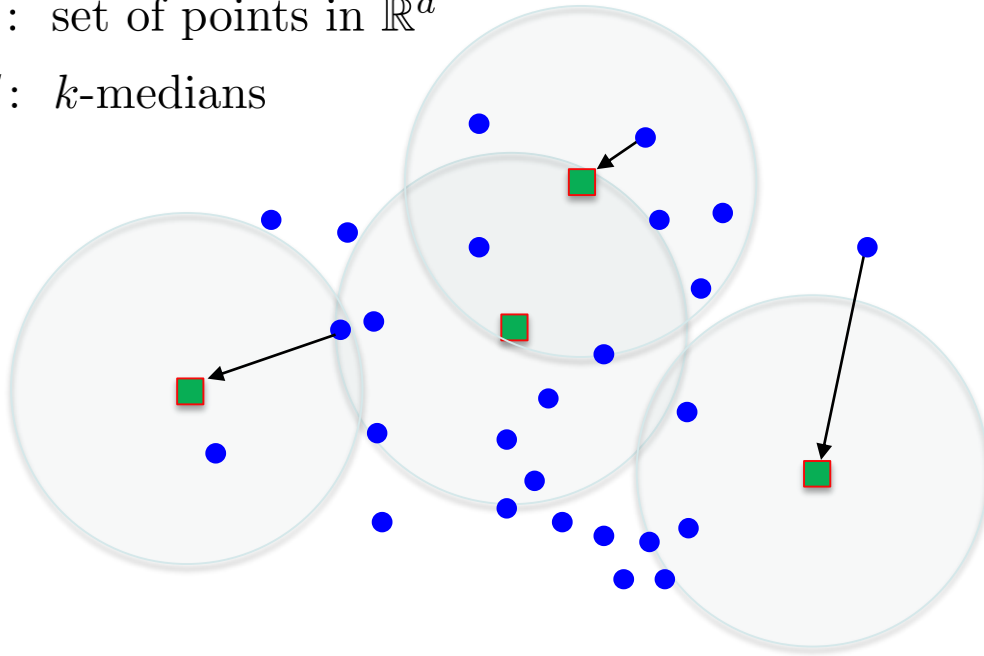
APPLICATION

ϵ -approximations

CLUSTERING

P : set of points in \mathbb{R}^d

C : k -medians



Goal: compute a small-sized set A such that

$$\text{for any } C, |C| = k: \sum_{p \in A} \text{distance}(p, C) \cdot w(p) = (1 \pm \epsilon) \cdot \sum_{p \in P} \text{distance}(p, C)$$

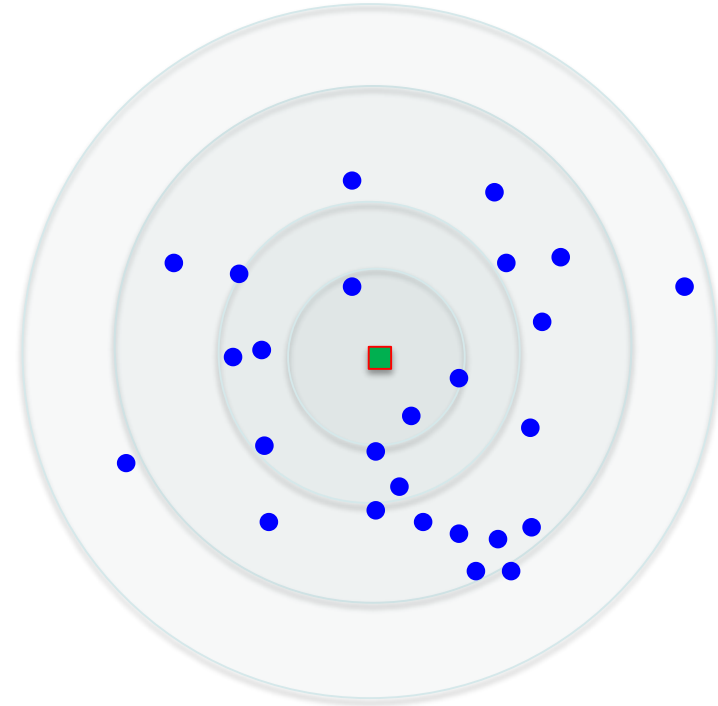
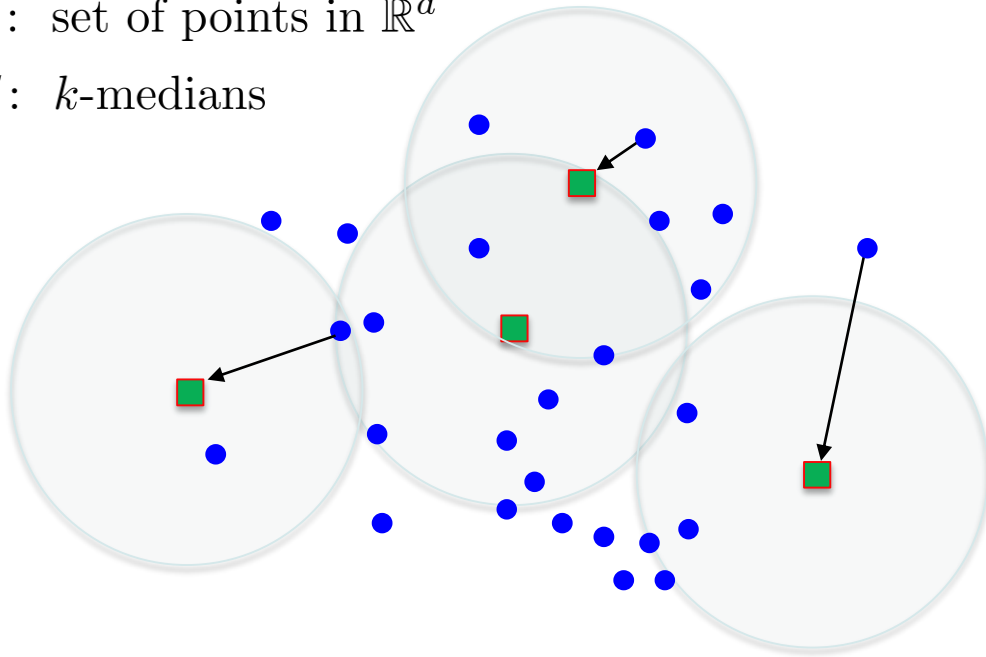
[Feldman, Langberg 2011]

\implies computing optimal clustering on A gives approximate clustering on P

CLUSTERING

P : set of points in \mathbb{R}^d

C : k -medians



Theorem : there exists assignment of weights, $w: P \rightarrow \mathbb{R}$, such that

any ϵ -approximation A for union of k balls w.r.t. these weights satisfies

$$\text{for any } C, |C| = k: \quad \sum_{p \in A} \text{distance}(p, C) \cdot w(p) = (1 \pm \epsilon) \cdot \sum_{p \in P} \text{distance}(p, C)$$

[Feldman, Langberg 2011]

\implies computing optimal clustering on A gives approximate clustering on P

CONCLUSION

random sampling on combinatorial structures

structures improve sampling bounds and analysis

probability, statistics, learning

‘altering’ samples gives better bounds

newer algorithms not much more complicated

combinatorics, geometry

applications of these compact structures

optimisation, graphs, algorithms

step 1. show the existence of a probability distribution on a set system

step 2. computing approximation under this distribution

LP, FL

Thank you