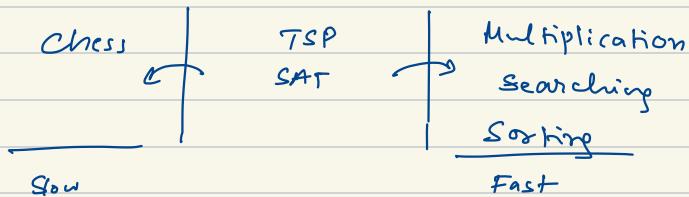


## Computational Hardness

Given a set of locations, find the shortest possible route to visit all locations by visiting each exactly once. (TSP)



$P = \{ \text{Problems solvable in polynomial time by deterministic TM} \}$

$n^{O(1)}$  time  $\rightarrow$  polynomial

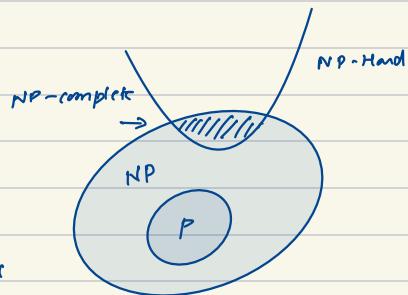
$NP = \left\{ \begin{array}{l} \text{Decision problems solvable in polynomial time} \\ \text{by non-deterministic TM} \end{array} \right\}$

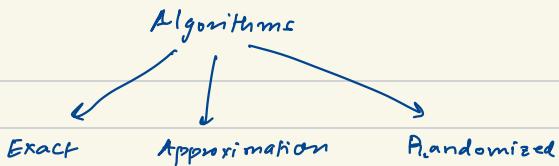
I/P  $\rightarrow$  [Verifier]  $\rightarrow$  Yes

A problem  $X$  is

- NP Complete if  $X \in NP \wedge X$  is NP-Hard
- NP Hard - A class of problems already established.

If we can show a reduction to these then we call such problems NP Hard





## Course Content

### First Half

Approximation Algorithms : Ind-set, Set Cover, Hitting Set,  
TSP, Clustering (K-means, k-median),  
Steiner Tree

Parametrized Algorithms : Kernelization, Color Coding, Parametrized approximation,  
FPT Approximation, Lossy Kernelization, Bi-dimensionality

### Second Half

Dynamic Aspects : Data is not static

Paradigms:

- Online } → Reveal the input piece by piece
- { - Dynamic → stochastic Model where input is unknown
- Streaming      but is drawn from some known distribution

Query complexity, update time of the data structure

Graphs : Connectivity, reachability, apx. distance oracles

Geometry: LSH, point location, range searching.

Other: Succint DS, External memory, sketching

## Approximation Algorithm

Design an algorithm that strictly runs in polynomial time ( $n^{O(1)}$ )

Output is allowed to be a "provable" factor away from the optimal sol?

### Maximization Problems

Eg: Ind. Set

Variable  $\alpha \geq 1$ .  $\alpha$ -approximation

if we output a sol? that is  
 $(\frac{1}{\alpha})$ -factor to the output

### Minimization Problems

Eg: Hamiltonian Cycle.

$\alpha$ -appx if we output a sol?  
that is at most  $\alpha \times$  opt.

## Polynomial-Time Approximation Scheme

An algorithm (given with some parameter  $\varepsilon > 0$ )

for any input, output a sol? within a factor  $(1 + \varepsilon)$  of the optimal sol?  
that runs in  $n^{f(\varepsilon)}$ .

For some computable function  $f$ .

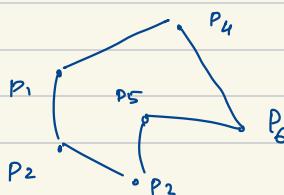
Approx. Algos - Vazirani, Shmoy's

(Running time is polynomial on  $n$ , some  $\varepsilon$ )

Williamson

## Travelling Salesman Problem (TSP) ~1930s

Given a list of cities ( $P \subseteq \mathbb{R}^2$ ) and distances between each pair of cities,  
goal is to compute the shortest possible route that visits each city  
exactly once.



Decision version: Given length  $L$ , Is it possible to find a solution of length atmost  $L$ .

Graph I/p:  $G(V, E, W)$

O/p - Visit all vertices without repetition

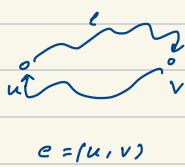
Minimizing the sum of edge weights

### \* Hamiltonian Cycle Problem

→ Hamiltonian Cycle is NP-Complete [Richard Karp 70's]

→ In fact, No constant factor approx is possible

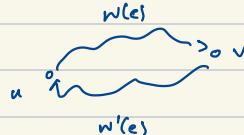
Symmetric



$$e = (u, v)$$

$$\underline{w(e)}.$$

Asymmetric



$$w(e)$$

### # Metric TSP

- $d(x, y) \geq 0$
- $d(x, y) = d(y, x)$
- Triangle Inequality:  $d(x, y) + d(y, z) \geq d(x, z)$



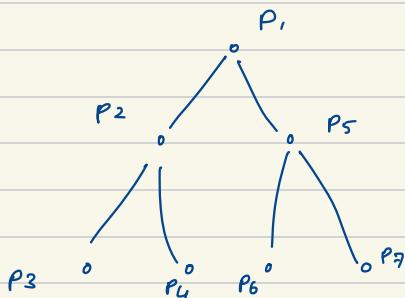
$S$  - Set of Edges

$\text{Cost}(S)$  - Sum of the weights of edges (min)

(Same as finding minimum weight hamiltonian Path/ Cycle.)

Base structure:

- Min. Spanning Tree (Kruskal)



Repeat not allowed

P<sub>1</sub> P<sub>2</sub> P<sub>3</sub> P<sub>2</sub> P<sub>4</sub> P<sub>1</sub> P<sub>5</sub> P<sub>6</sub> P<sub>5</sub> P<sub>7</sub> P<sub>5</sub> P<sub>1</sub>

↓  
This is not a valid tour

→ We know that there is no constant factor approximation.  
And we haven't really done anything special

→ Adding new paths to remove re-visitation will decrease cost  
Because of the triangle inequality in a metric space.

- DFS Traversal

- Delete the duplicates

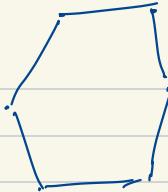
↳ Due to triangle inequality, Duplication is always possible

# Analysis: Obviously, algorithm is polynomial time.

→ Every edge of the MST is travelled twice (pre duplication)

$$\text{cost}(C) \leq 2 \times \text{cost(MST)}$$

$$\text{Cost}(\text{MST}) \leq \text{cost}(\text{opt}) \rightarrow$$



Cycle  $\rightarrow$  can delete the longest edge to get  
 $\text{SpanningTree} \geq \text{MST}$ .

$$\therefore \text{Cost}(C) \leq 2 \times \text{Cost}(\text{MST}) \leq 2 \times \text{Cost}(\text{opt})$$

2-factor approximation!

Q Can we do better?

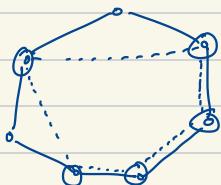
### PROPERTY 1

$$\text{I/p : } G(v, E) \longrightarrow \text{OPT}(G)$$

Take a subset  $S \subseteq V$

Induced subgraph  $G[S]$   $\xrightarrow{\text{OPT}(S)}$

$$\text{OPT}_S \leq \text{OPT}_G$$



### PROPERTY 2 :

Perfect Matching (Can be computed in Polynomial Time )

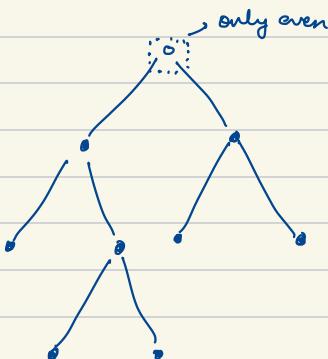


Minimum cost perfect matching is a perfect matching with smallest cost.  
 (Poly.)

### PROPERTY 3 :

Eulerian tour (Circuit) : Start from one vertex and return back to it after visiting all edges only once. (Allowed to repeat an edge but not repeat any vertices)

If graph has even degree vertices then you always have an Eulerian Cycle.



Compute an Eulerian circuit

→ if we match odd degree vertices, their degree increases by 1

→ However, we can have an odd number of odd degree vertices, can't we?

≤ How many odd degree vertices can we have?

$$\sum_{v \in V} d(v) = 2 \times |E|$$

Since this is even, ignoring twice even degree vertices, we must have an even number of odd degree vertices.  
Every edge is counted

Now we have a graph with all even edges.

- Compute Eulerian circuit
- Delete duplication

Analysis:

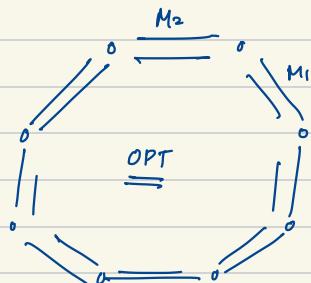
$$\text{Cost}(C) = \text{cost}(\text{MST}) + \text{cost}(\text{Matching})$$

↓ removing duplication

$$\text{Cost}(C')$$

$$(\text{cost}(C') \leq \text{cost}(C))$$

$$\begin{aligned}\text{cost}(C') &\leq \text{cost}(\text{MST}) + \text{cost}(\text{Matching}) \\ &\leq \text{cost}(\text{opt})\end{aligned}$$



$$\text{cost}(M_1) \leq \text{cost}(\text{opt})$$

$$\text{cost}(M_2) \leq \text{cost}(\text{opt})$$

$$\text{cost}(\text{opt}) \geq \frac{\text{cost}(M_1) + \text{cost}(M_2)}{2}$$

$$\text{cost}(M) \leq \frac{\text{cost}(M_1) + \text{cost}(M_2)}{2} \leq \frac{1}{2} \text{cost}(\text{opt})$$

$$\text{cost}(C') \leq \text{cost}(\text{opt}) + \frac{\text{cost}(\text{opt})}{2} = 1.5 \text{cost}(\text{opt})$$

Metric TSP

- 2 - apx. (MST doubling)

- 1.5 - apx. (Christofides Algo.) [1976]

$1.5 - \varepsilon$  apx.  $\varepsilon \sim 10^{-20}$

→ No apx. is possible if the distance is arbitrary

Q Does there exist a PTA for Metric TSP

$(1+\epsilon)$ -APX  
 $n^{O(f(\epsilon))}$  time

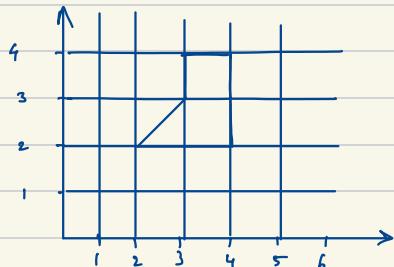
Thm: There can't be a PTAS  $2^{20/21\epsilon}$ -apx unless  $P=NP$

Restrict the metric: Euclidean Metric

I/p: A set of points in  $\mathbb{R}^2$ , with Euclidean distance

$$d(x,y) = \|x-y\|_2$$

O/p: Find the shortest route that visits all pb.



It is not known whether the problem is in NP.  
It is however known that the problem is NP Hard.

### Sum of Square Roots

Given a set of positive integers  $a_1, a_2, \dots, a_k, +$

Decide

$$\sum_{i=1}^k a_i \leq t$$

$$\{a_1, a_2, \dots, a_k\} \quad \{b_1, b_2, \dots, b_k\}$$

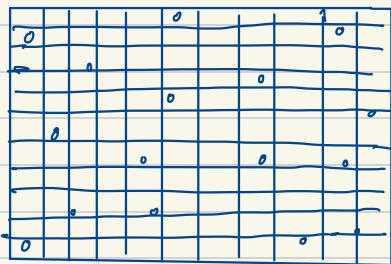
$$\sum_{i=1}^k a_i \leq \sum_{i=1}^k b_i$$

Affects all Euclidean Problems

- Euclidean TSP
- Euclidean MST
- Steiner Tree
- Matching

## PTAS for Euclidean TSP [Arora (1998) & Joe Mitchell (1999)]

- Rounding the instance
- Partitioning: Exploit the structure of the instance by breaking it into "nice" instances.
- Apply dp to the instance.



$R_2$

① Wlog you snap each point to the closest grid point

We must first prove that doing so doesn't lead to much "loss"

$\epsilon$  - "nice" instance

Def<sup>2</sup> An instance of Euclidean TSP is  $\epsilon$ -nice if :

1. Every point has integral coordinates in the interval  $[0, \epsilon(\frac{n}{\epsilon})^2]$
2. Any 2 different points have distance atleast 4.

→ Take a small bounding box (axis-parallel)

→ Longer side = L

→ Translate the instance & root at origin, scale  $L = \lceil \frac{8n}{\epsilon} \rceil$

$\left. \right\} \text{=}'$

Lemma - I is i/p  $\Rightarrow \text{OPT}_I$  is optimal tour

$I'$  is  $\epsilon$ -nice instance  $\Rightarrow \text{OPT}_{I'}$  is optimal tour

$$\text{OPT}_{I'} \leq (1+\epsilon) \text{OPT}_I$$

We know that  $\text{OPT}$  is at least  $\epsilon L$ ,

Draw a fine grid with spacing  $\frac{\epsilon \times L}{2n}$

- Map every point to its closest grid pt.  
(multiple pts. could be mapped to same)
- All points have int coords

$$k = \lceil \frac{8n}{\epsilon} \rceil \in o(n/\epsilon)$$

$$\frac{\epsilon L}{2n} \geq \frac{\epsilon L}{2n} \times \frac{8n}{3} = 4$$

→ Mapping each pt. in  $I$  has moved  $\frac{\epsilon L}{2n}$   
(Every edge in the set  $I$  changed  
at most  $\frac{\epsilon L}{2n}$ )

Cost:  $\text{OPT}_I + \epsilon \times L$

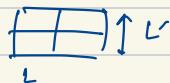
$$\leq \int_L^U L \leq \text{OPT}_I$$

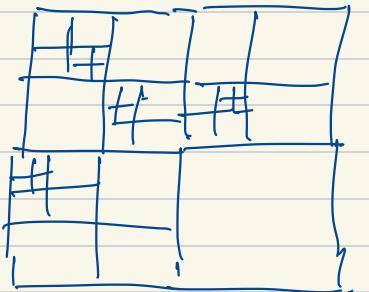
$$\text{OPT}_I + \epsilon \times \text{OPT}_I$$

$$\text{OPT}_{I+1} \leq (1+\epsilon) \text{OPT}_I$$

## ② Partition the Space

- Extend the bounding box to a square with new side length  $L'$   
 $L'$  is the smallest power of 2



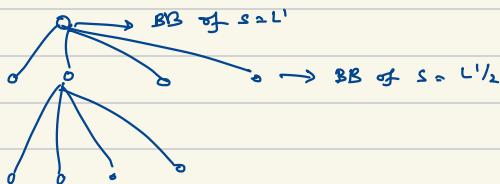


Quadtrees

- ★ Recursively partition the box/square into four equal sized squares until the side length becomes " $\epsilon$ "  
 $\downarrow$       ( $L'$  is a power of 2)

At the end of this partition, we can guarantee that there will be at most 1 pt. in each square  
 (since distance b/w 2 pts is atleast 4)

- Each pt. is separated
- One pt. in each "non-empty" square



height:  $\log(L')$

Partitioning terminates after  $O(\log L')$  steps

$$\text{Height of quadtree} = O(\log L') = O(\log(n/\epsilon))$$

Idea: Apply dynamic programming to the Quadtree

↓  
Solve for each square that are leaves

↓  
Bottom-up combine.

1/8

### Euclidean TSP

Given a set of  $n$  points in  $\mathbb{R}^d$  with distances  $\forall x, y \quad d(x, y) = \|x - y\|_2$

Output: shortest tour that visit all points

PTAS ( $1 + \epsilon$ ) in  $n^{O(\ell(\frac{1}{\epsilon}))}$

-  $\epsilon$ -niceness

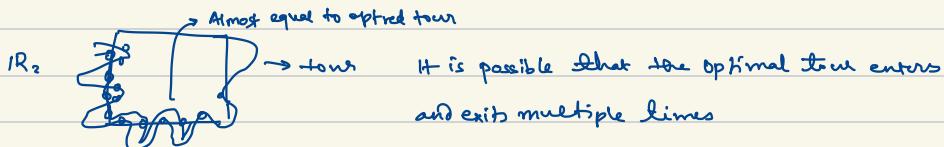
- quadtree partitioning

- DP

$\epsilon$ -niceness - 1. Integral

2.  $d(x, y) \geq 4$

Quadtree - Every cell must have exactly 2 point



$(1 + \epsilon)$ -opt → Too many calls to the box will lead to exponential blowup.

→ Add "portals" on the box

→ Show that optimal tour must only enter and exist through these portals

→ Show that entry & exits through these portals is bounded.

- Limits the # of interactions.

# of portals

- Accuracy
- Running time

Trade off

Select m to be a power of 2  $m \in \left[ \frac{k}{\varepsilon}, \frac{2k}{\varepsilon} \right]$

For each square,

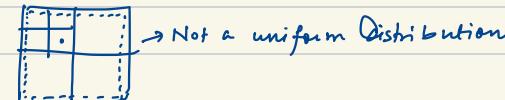
- Put portals in corners
- Put  $(m-1)$  portals equally spaced
- Few many portals

### Portal Respecting Tours

Def: p-tour enters & exits through portals

- Length of p-tour  $\leq (1 + \varepsilon)$  -optimal

- Detours can add much more cost.



Sol<sup>1</sup>: Randomize

1. Translate the grid by a "random offset" at most  $1/\varepsilon$  in each coordinate
2. Points remain grid points
3. With high probability, the points are "nicely" concentrated.
4. Higher-levels in the partition (quadtree) have more portals  $\rightarrow$  have fine grained tour

Defn : (a,b)-dissection : Origin of the grid is translated by  $(-a, -b)$

Thm:  $(a,b)$  picked up uniformly @ random  $[0, \frac{1}{2}]$   
with prob atleast  $\frac{1}{2}$

p-tour such that

$$\text{cost(p-tour)} \leq (1+4\epsilon) \times \text{opt}$$

Prof :- Extend non p-tour to a p-tour.

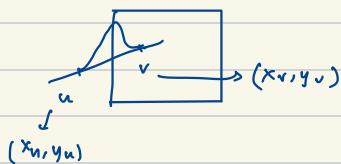
\* For each vertical/horizontal line  $\ell$

$$t(\ell) = \# \text{ times p-tour intersects } \ell$$

$$T = \sum_{\ell} t(\ell)$$

claim:  $T \leq 2 \times \text{opt}$

- e crosses  $x+1$  vertical lines
- e crosses  $y+1$  horizontal lines



e crosses  $(x+y+2)$  lines

$$\sqrt{2(a^2+b^2)} \geq (a+b) \quad \text{--- (1)}$$

$$\forall x,y \quad d(x,y) \geq 4 \quad \text{--- (2)}$$

$$\begin{aligned} (x+y+2) &\leq \sqrt{2(x^2+y^2)} + 2 \\ &\leq 2 \underbrace{\sqrt{x^2+y^2}}_{\text{opt}} \end{aligned}$$

## Bound - expected length of the detour

Detour might cross

$$|x_u - x_v| + |y_u - y_v| + 2$$

$i$  of the quad-tree

$$\frac{L'}{2^{im}}$$

$$\rightarrow \text{if } l \text{ is in level } i \leq \frac{L'}{2^{im}}$$

Q :- What is the probability that after "randomshift"  $l$  crosses  $e$  at level  $i$ ?

$\rightarrow l$  could be mapped to  $L'/2$  many lines [translated by  $(0, L'/2)$ ]

$\rightarrow 2^{i-1}$  many lines of level  $i$

$$\frac{2^{i-1}}{L'/2} = \frac{2^i}{L'}$$

ub on expected length

$$\sum_i^K \frac{2^i}{L'} \times \frac{L'}{2^{im}} \leq \epsilon$$

By linearity of expectation

$2\epsilon \text{OPT}$

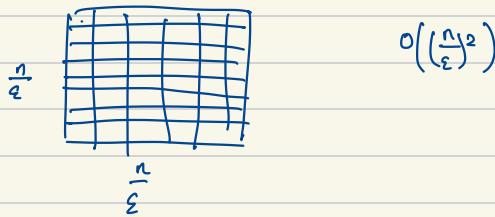
Markov Inequality

$$\Pr[\text{total length increase of detours} \geq 4\epsilon \text{ OPT}] \leq \frac{2\epsilon \text{ OPT}}{4\epsilon \text{ OPT}} = \frac{1}{2}$$

## De randomization

- fixed  $\epsilon$

grid shifting by trying all possibilities



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

## Final Step (DP)

Given  $(a_{ib})$  - dissection, get p-tour

Introduces state -

- a square

- any set of possible ways of entering/exiting  
the square.

# states must be polynomial

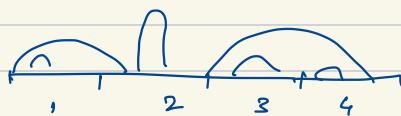
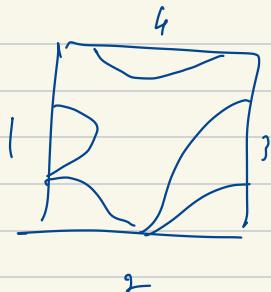
$$1 + 4 + 4^2 + \dots + L^2 = O\left(\frac{n^2}{\epsilon^2}\right)$$

Lemma: w.l.o.g. a p-tour is well-behaved 2-light



- 4m portals
- use one portal  $\{0, 1, 2\}$  times

$$\begin{aligned} 3^{4m} &\rightarrow m = o(k/\varepsilon) \\ &= o(\log n/\varepsilon) \\ &= O((\log n/\varepsilon)) \end{aligned}$$



( ) ( ) ( ) ( )

$$\begin{aligned} r^{\text{th}} \text{ Catalan Number} \rightarrow C_r &= \frac{(2r)!}{r+r!} = o(2^{2r}) \\ &= o(2^{8m}) \quad m \text{ is } o(k/\varepsilon) \end{aligned}$$

- Translate them into paths
- Discard anything that intersects

$$m = O\left(\log\left(\frac{n}{\epsilon}\right)/\epsilon\right) \Rightarrow \text{Entry/Exit} = O(n/\epsilon)$$

Computation of values:

$$A[(s_1, t_1), (s_2, t_2), \dots (s_r, t_r)]$$

↳ Compute the whole table

## CLUSTERING

→ learning, searching, data-mining, ...

→ Data represented as points in  $\mathbb{R}^d$

→ General metric space  $(X, d)$

$d$  → distance  $d: X \times X \rightarrow [0, \infty)$   
 Space metric  
 $(\mathbb{R}^d)$

(set)

$(X, d)$  is a metric if it satisfies -

(i) if  $x=y$   $d_X(x, y)=0$

(ii)  $\forall x, y$   $d_X(x, y) = d_X(y, x)$

(iii)  $\forall x, y, z$   $d_X(x, y) + d_X(y, z) \geq d_X(x, z)$

Assumption: Given  $x, y$   $d_X(x, y)$  can be computed in  $O(1)$  time.

Norm

→ norm defines distances between points  
 $p, q \in \mathbb{R}^d$

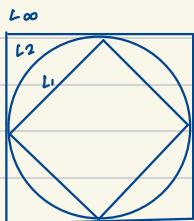
$$\|p - q\|_t = \left( \sum_{i=1}^d |p_i - q_i|^t \right)^{\frac{1}{t}} \quad \text{for } t \geq 1$$

$t=2$  : Euclidean Norm

$t=1$  :  $L_1$ -norm, Manhattan norm

$$L_\infty\text{-norm} : \|p - q\|_\infty = \lim_{t \rightarrow \infty} \left( \sum |p_i - q_i|^t \right)^{\frac{1}{t}}$$
$$= \max_i \{|p_i - q_i|\}$$

Triangle inequality holds for  $L_\infty$  too, it's called Minkowski Inequality



For any pt.  $p \in \mathbb{R}^d$

$$\|p\|_t \leq \|p\|_s \quad \text{if } t > s$$

Lemma - For any  $p \in \mathbb{R}^d$

$$\frac{\|p\|_1}{\sqrt{d}} \leq \|p\|_2 \leq \|p\|_1$$

Proof :-  $p = (p_1, p_2, \dots, p_d) \Rightarrow p_i \geq 0 \forall i$

Const.  $\alpha$

obj:  $f(x) = x^2 + (\alpha - x)^2$  minimized if  $x = \alpha/2$

$$\text{Set } \tau = \|p\|_1 = \sum_{i=1}^d |p_i|$$

By symmetry & obs. on  $f(n)$

$$\|P\|_2 \geq \sqrt{d(\alpha/\alpha)^2} = \frac{\|P\|_1}{\sqrt{d}}$$

## K-centre clustering

Metric Space  $(X, d)$

I/P : A set of points  $P$ ,  $|P|=n$

O/P : Find  $n$ -clustering (set of  $n$  clusters)

such that each pt. is assigned to its nearest centre

Set of clusters  $C$  prime cluster

$$\text{cluster}(c, \bar{c}) = \left\{ p \in P \mid d_M(p, \bar{c}) = d_M(p, c) \right\}$$

max dist in the set  
~~~~~  
set.

$n$  dim pt.  $p_i = (d(p_{i,1}, c), d(p_{i,2}, c), \dots, d(p_{i,n}, c))$

ith coordinate  $d(p_i, c)$  dist. to pi  
to its closest centre.

---

O/P :- Find a set of  $k$ -centres  $C \subseteq P$  such that the maximum distance of  $n$  pt. in  $P$  to its closest centre is minimized.  
(Facility Location)

## Formal Statement

Given a set of  $k$ -centres  $C$ ,

$$\|P\|_{loc} = \max_{p \in P} d(p, c)$$

Find  $c_r$  s.t.  $\|Pc\|_\infty$  is minimized

$$\text{opt}_\infty(P, k) = \min_{C \subseteq P, |C|=k} \|Pc\|_\infty$$

→  $C_{opt}$

→ NP-Hard

→ Hard to approximate within 1.86

→ 2-approximation in the Euclidean Case in  $\mathbb{R}^2$

### Greedy-Algo

- Start by picking an arbitrary point,  $\bar{c}_1 : C_1 = \{\bar{c}_1\}$

- Compute the distances for each  $p \in P$ , from  $\bar{c}_1$

- Take pt. with worst dist.

$$(r_1 = \max_{p \in P} d_1(p))$$

Say  $\bar{c}_2$

add  $c_2 : C_2 = C_1 \cup \{\bar{c}_2\}$

i<sup>th</sup> step:  $C_i = C_{i-1} \cup \{\bar{c}_i\}$

We have atmost  $k$ -centres  $\Rightarrow$  atmost  $k$  iterations to the algorithm.  $O(n)$  time to compute distances & update.

$O(nk)$  time

$O(n)$  space

For each  $p \in P$ , a single variable  $d[p]$  with its current list to the closest pair then

$$d[p] \leftarrow \min(d[p], d_M(p, \bar{c}_i))$$

### Def<sup>n</sup>

A ball of radius  $r$  around a pt.  $p \in P$  is a set of points in  $P$ , with dist. atmost  $r$  from  $P$

$$b(p, r) = \{q \in P \mid d_M(p, q) \leq r\}$$

Remark -  $K$ -centre is essentially covering  $k$  with  $k$  balls of min Radii

Thm. Greedy algo computes a set  $K$  of  $k$ -centres s.t.  $K$  is 2-approx.  
 $\|PK\|_\infty \leq 2 \times \text{opt}_{\infty}$   
 and takes  $O(nk)$  time.

Proof :- Running Time ✓

→ By Def<sup>n</sup>  $\pi_K = \|PK\|_\infty$

Let  $\bar{c}_{K+1}$  is the pt. realizing  $r_K = \max_{p \in P} d(p, K)$

$$C = K \cup \{\bar{c}_{K+1}\}$$

By the definition of  $\pi_i$ ,

$$r_1 \geq r_2 \geq \dots \geq r_K.$$

$$i < j < K+1$$

$$d_M(\bar{c}_i, \bar{c}_j) \geq d_M(\bar{c}_j, c_{j-1}) \quad r_{j-1} \geq r_K$$

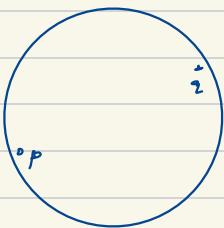
- the distance between any pair of points in  $C_i$  is at least  $r_k$

Opt :- Clusters  $P$  by using  $k$  balls

By Triangle inequality - any 2 points within such a ball  
are with a dist. atmost  $2 \times \text{opt}$



None of the balls contain 2 points for  $C \subseteq P$



### Greedy Permutation

Let this run till we exhaust all points

$$\langle p \rangle = \langle \bar{G}_1, \bar{G}_2, \dots, \bar{G}_n \rangle$$



$$\langle r_1, r_2, \dots, r_n \rangle$$

Def<sup>n</sup>  $r$ -packing . A set  $S \subseteq P$ :

(i) Covering Property: All points in  $P$  are within distance  
of atmost  $r$  from  $s$ .

(ii) Separation Property:  $d_{\mu}(p, q) \geq r$

\*  $r$ -packing gives compact representation

\* Greedy permutation gives such a representation

then

$$\langle \bar{c}_1, \bar{c}_2, \dots, \bar{c}_n \rangle, \langle r_1, r_2, \dots, r_n \rangle$$

for any  $i$ , we have  $c_i = \langle \bar{c}_1, \dots, \bar{c}_i \rangle$   
is an  $n$ -packing of  $P$

Proof :- By contradiction,

$$r_{k-1} = d(\bar{c}_k, c_{k-1}) \quad \forall k=2, \dots, n$$

For  $1 \leq k \leq n$

$$d_M(\bar{c}_i, \bar{c}_k) = r_{k-1} \geq r_i$$

### K-medians cluster

A set  $P \subseteq X$ ,  $|P|=n$ . A parameter  $K$ .

Find a set of  $K$  points  $C \subseteq P$ , s.t. the sum of distances of the pts. in  $P$  to its closest centre is minimized.

### Clustering prices

$$\|p_C\|_1 = \sum_{p \in P} d(p, c)$$

### Objective fn

$$OPT(P, K) = \min_{\substack{C \subseteq P, \\ |C|=K}} \|p_C\|_1$$

### Optimal Set of Centre - $C^{OPT}$

Local Search: Move  $\text{sol}^0$  to  $\text{sol}^1$  in the space of the candidates  $S\Omega^n$  (the search space) by applying local change

Continue until - end up on optimal

or

we exhaust the running time

For this problem  $n^k$  possible solutions - so an exhaustive local search should lead to the optimal  $\text{sol}^n$ .

### Notations

A set  $U = \{p_c \mid c \in \mathbb{P}^k\}$

$$\text{OPT}_{\infty}(P, k) = \min_{Z \in U} \|Z\|_{\infty}$$

$k$ -centre

$$\text{OPT}_1(P, k) = \min_{Z \in U} \|Z\|_1$$

$k$ -median

Claim: For any set  $P$ ,  $|P| = n, k$

$$\text{OPT}_{\infty}(P, k) \leq \text{OPT}_1(P, k) \leq \sum_{i=1}^k \text{OPT}_{\infty}(P_i)$$

Proof...  $p \in \mathbb{R}^n$

$$\|P\|_{\infty} = \max_{i=1}^n |p_i| \leq \sum_{i=1}^n \|p_i\|_1 = \|P\|_1$$

$$\|P\|_1 = \sum_{i=1}^n \|p_i\|_1 \leq \sum_{i=1}^n \left( \max_{j=1}^n \|p_j\|_1 \right) \leq n \max_{i=1}^n |p_i| = \sum_{i=1}^k \text{OPT}_{\infty}(P_i)$$

$C$ -set of centres;  $|C| = k$

realizing  $\text{opt}_1(p, k)$  i.e.  $\text{opt}_1(p, k) = \|p\|_1$

$$\text{opt}_{\infty}(p, k) \leq \|p\|_{\infty} \leq \|p\|_1 = \text{opt}_1(p, k)$$

Similarly,

$$k \text{ realizing } \text{opt}_{\infty}(p, k)$$

$$\begin{aligned}\text{opt}_1(p, k) &= \|p\|_1 \leq \|p\|_1 \\ &\leq n \times \|p\|_{\infty} \\ &= n \times \text{opt}_{\infty}(p, k)\end{aligned}$$

$\rightarrow 2d$ -factor for  $k$ -median

This is not good. We can't effectively modify the Greedy Algorithm either since it starts getting very inefficient.

Here we need Local Search. Given an initial algorithm, we try to improve it using Local Search.

We use this greedy algorithm as the first step of Local search.

#  $L$  is  $2d$ -approximation

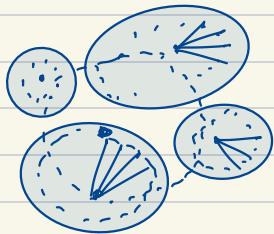
Improve - parameter  $0 < \delta < 1$

$w \in [n]$   
↳ curr

## Local Search -

→ Set :  $L_{curr} \leftarrow L$

→ At each iteration



This is very bad → non-centres  
We try and kick out a point &  
add a new pt. to improve cost  
by some factor  $\delta$

$$\text{swap} \rightarrow K \leftarrow (L_{curr} \setminus \{\epsilon\}) \cup \{ \underset{\hookrightarrow \text{new}}{\delta} \}$$

$$\text{if } \|P_k\|_1 \leq (1-\delta) \|P_{curr}\|_1$$

→ Continue "swap" as long as it satisfies the constraint

→ Return  $L_{curr}$

## Running Time

An iteration takes  $O(n \times k)$  swaps  $\xrightarrow{\text{naively}}$   
 $n-k$  can be swapped in  $k$  to be swapped in

$$\text{Since } \frac{1}{1-\delta} \geq 1+\delta$$

$$O\left((nk)^2 \log \left( \frac{\|P_L\|_1}{(1-\delta)} \frac{\|P_L\|_1}{opt_1} \right) \right) = O\left((nk)^2 \log_{1+\delta} n \right)$$

$$= O\left(\frac{(nk)^2 \log n}{\delta}\right)$$

## K-means

Set  $P \subseteq X, K,$

Find  $K$  points  $C \subseteq P, |C| = k,$

$$\|P_C\|_2^2 = \sum_{p \in P} (d_H(p, c))^2$$

Obj:-

s.t.  $\|P_C\|_2^2$  is minimized

$$\text{opt}_2(P, k) = \min_{C, |C|=k} \|P_C\|_2^2$$

$O(1)$  - factor for k-means too.

Same strategy, apply standard greedy & then Local Search

Then  $0 < \varepsilon < 1$

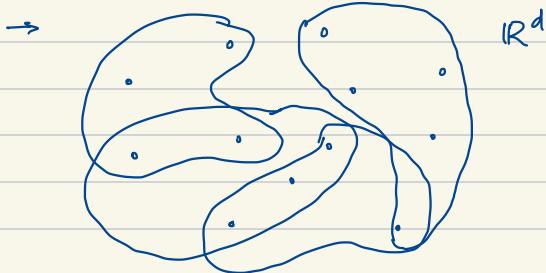
$$(25 + \varepsilon) \text{ factor app} \quad O(n^2 k^3 \frac{\log n}{\varepsilon})$$

## VC-Dimension

→ range space  $(X, \mathcal{R}) = \mathcal{S}$

$X$ : ground set (finite / infinite)

$\mathcal{R}$ : family of subsets of  $X$



- Sampling for exit polls for example
- blue blobs are groups.
- possible for  $e$  to be in multiple groups

→ Consider finite subset of  $X$  as the estimating ground set.

$X \setminus N$

Def<sup>n</sup> (Measure)

$X$  fixed subset of  $\mathbb{R}^d$ . For a range  $\mathcal{R} \in \mathcal{S}$ ,  
measure for  $\mathcal{R}$

$$\bar{m}(x) = \frac{|\mathcal{R} \cap x|}{|x|}$$

Estimate

For a subset  $N$  (multi-set) of  $X$ , its estimate of the measure of  $\bar{m}(x)$ , for  $x \in \mathcal{R}$

$$\bar{\delta}(x) = \frac{|\mathcal{R} \cap N|}{|N|}$$

Q Can we get methods to generate  $N$  st

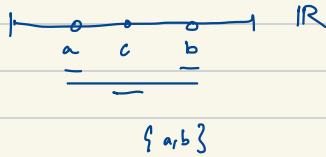
$$\hat{s}(x) = \hat{m}(x), \forall x \in \mathbb{R}$$

VC-Dimension  
↳ chernomskis  
Vapnik (1971)

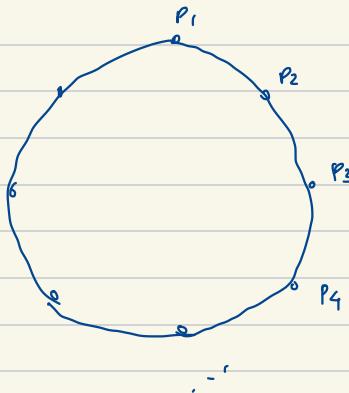
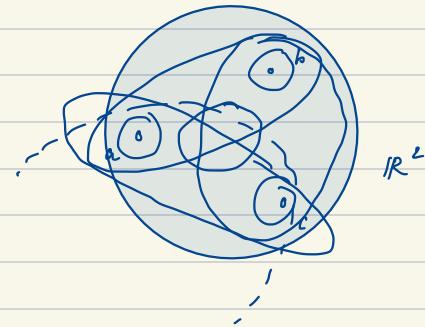
Def<sup>n</sup>  $S = (X, R)$  for  $Y \subseteq X$

$R_Y = \{s \cap Y \mid s \in R\}$  ↪ if  $2^{|Y|}$  then we call it shattered by  $\mathbb{R}$   
be the projection of  $R$  on  $Y$

The orange space  $S$  is projected to  $S_{1Y} = (Y, R_{1Y})$



can't shatter



$$\dim_{VC}(S) = \infty$$

## Complement

$$S = (X, R) \quad , \quad s = \dim_{VC}(S)$$

$$\bar{S} = (X, \bar{R}) \quad , \quad \text{where } \bar{R} = \{x^{\forall r} \mid r \in R\}$$

Q What is the VC-dimension of  $\bar{S}$ ?

A subset  $B \subseteq X$  is shattered in  $\bar{S}$ .

iff it is shattered in  $S$ .

$$\text{For any } z \in B, (B \setminus z) \in R|_B \Rightarrow z = B \setminus (B \setminus z) \in \bar{R}$$

$$\text{Thm: } \dim_{VC}(S) = \dim_{VC}(\bar{S})$$

## Local Search

- Let  $X$  be the set of arbitrary subset of  $k$ -centres (greedy sol<sup>n</sup> k-centres)

while True :

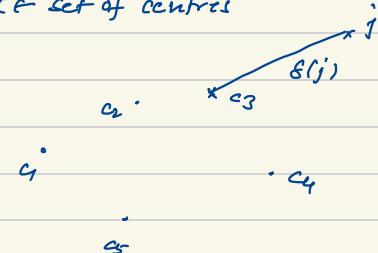
swap : if  $i \in X$  and  $i' \in F \setminus X$   
 $\text{cost}(x - i + i') < \text{cost}(x)$  :  
 $x \leftarrow x - i + i'$

otherwise :

break;

Claim: 5-approx

$X \in$  set of centres



$\text{opt} = X^* - \text{optimal set of } k\text{-centre}$

Nearst

$i' \in X^*$

$i \in X$

→ For each centre chosen in  $X$ , the nearest centre in  $X^*$  is what we want

$\min(i, i')$

Inverse

For the centers in  $X^*$ , inverse to the nearest map.

Bijection

Claim:- For any  $j \in X$ ,

$$d(\text{nearest}(x^*(j)), j) \leq d_j + 2d_j^*$$

## Half Space

If  $d$  dimensional then  $d+1$  dimensional Ray spaces

Let  $R$  be the set of closed half spaces in  $\mathbb{R}^d$ .

Claim :-

$$P = \{ p_1, \dots, p_{d+2} \} \quad \text{a set of points in } \mathbb{R}^d$$

Real nos.  $\beta_1, \beta_2, \dots, \beta_{d+2}$  (not all are zero)

$$\text{s.t.} \quad \sum_i \beta_i p_i = 0 \quad \& \quad \sum_i \beta_i = 0$$

Proof :-  $q_i = (p_i, 1)$  for  $i=1, \dots, d+2$

pts are linearly ~~dependent~~ independent

$$(q_1, q_2, \dots, q_{d+2} \in \mathbb{R}^{d+1})$$

These are coefficients  $\beta_1, \dots, \beta_{d+2}$

$$\text{s.t.} \quad \sum_{i=1}^{d+2} \beta_i p_i = 0$$

Considering first  $d$ -coordinates of those pts.

implies  $\sum_{i=1}^{d+2} \beta_i p_i = 0$

$$\sum_{i=1}^{d+2} \beta_i p_i = 0$$

Similarly,  $(d+1)$  coordinates

$$\sum_{i=1}^{d+2} \beta_i = 0$$

Rander's Theorem :-  $p = \{p_1, \dots, p_{d+2}\}$

3 disjoint subsets  $C \& D \& P$

$$A.t. \quad CH(C) \wedge CH(D) = \emptyset$$

$$C \cup D = P$$

### Shattering dim

Property :- of range space  $(R)$  with  $VC\text{-Dim}(S)$   
means # of ranges grow polynomially on  $n$ ,  
(Generally exponential)

$$\text{Growth f} \ldots G_S(n) = \sum_{i=0}^{\delta} \binom{n}{i} \leq \sum_{i=0}^{\delta} \frac{n^i}{i!} \leq n^\delta \text{ for } \delta \geq 1$$

Sauer's Lemma :-  $S = (X, R)$

$$VC(S) = S, |X| = n, |R| \leq G_S(n)$$

Proof :-  $n=0, \delta=0$ , done

For some  $x \in X$   
element

$$R_x = \{ u \setminus \{x\} \mid u \subseteq X \in R \wedge x \setminus \{x\} \in R \}$$

and

contains

$$R \setminus n = \{ u \setminus \{x\} \mid u \in R \}$$
 doesn't contain  $n$

$$\underline{obs.} = |R_x| + |R \setminus n|$$

## Shatter Function

$S = (x, R)$ , Shatter function

$\Pi_S(m)$  is the maximum

# of sets that might be created  
by  $S$ , when restricted by the subsets of size  $m$

$$\Pi_S(m) = \max_{\substack{B \subseteq X \\ |B|=m}} |R|_B|$$

## Shattering Dim

↪ The smallest  $d$  s.t.  $\Pi_S(m) = O(2^d)$ ,  $\forall m$

Thm.  $S = (x, R)$  has shattering dimension  $d$ , then the  
VC-Dimension is bounded by  $O(d \log d)$

Proof:-  $N \subseteq X$  be the largest subset of  $X$  shattered  
by  $S$  &  $S$  is the cardinality

$$2^S = |R|_{N^S} \leq \Pi_S(|N|) \quad (S \geq \max(2, 2 \log c))$$

$$S \leq \log c + d \log S$$

$$\Rightarrow \frac{S - \log c}{\log S} \leq d$$

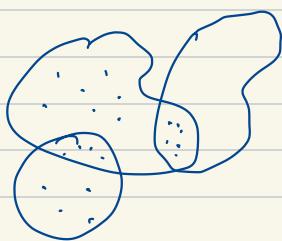
$$\frac{S}{2 \log S} \leq d \Rightarrow \frac{S}{\log S} \leq O(1 \times d)$$

$$f(x) = \frac{x}{\log x}$$

- mon. inc.  $x \geq e$
  - $f(x) \geq e^{-1}$  for  $x > 1$
  - $a \leq \sqrt{e}$ .  $f(a) \leq 4$  then  $x \leq a \log n$
- 

## $\varepsilon$ -net and $\varepsilon$ -Sampling

$\varepsilon$ -sample



$$S = (X, R)$$

$X$  is a finite subset of  $\mathcal{X}$

$0 \leq \varepsilon \leq 1$ , a subset  $C \subseteq X$  is an  $\varepsilon$ -sample for  $X$

if for any range  $R \in R$

$$|\bar{m}(r) - \bar{s}(r)| \leq \varepsilon$$

measure      ↓      estimate

(informally,  $\varepsilon$ -sample captures  $R$ , upto some " $\varepsilon$ " error)

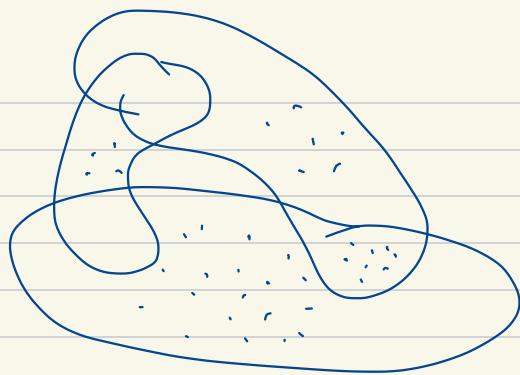
Theorem ( $\varepsilon$ -sample, VC Th): There is a tree constant  $c$ , s.t. if  $(X, R)$  is any range space with  $rc$ -dimension  $\delta$ .

$X \subseteq \mathcal{X}$  : finite subset of  $\mathcal{X}$  and  $\varepsilon, \phi > 0$ ,

a random subset  $C \subseteq X$  of cardinality

$$s = \frac{c}{\varepsilon^2} (\delta \log \delta / \varepsilon + \log 1/\phi)$$

with probability  $1 - \phi$



$\varepsilon$ -net :- A set  $N \subseteq X$  is an  $\varepsilon$ -net for  $X$  if for any range space  $(X, R)$  if  $m^*(r) \geq \varepsilon$

$$| \{ r \in X \mid \exists n \in N : r \in n \} | \geq \varepsilon |X|$$

then  $R$  contains atleast one pt. of  $N$  i.e.  
 $R \cap N \neq \emptyset$

### $\varepsilon$ -net theorem (HW8.7)

$S = (x, R)$  has VC-Dimension ( $C_S$ )  $\leq r$

$X$  is a finite subset of  $S$

Suppose  $0 \leq \varepsilon \leq 1$  &  $\delta < 1$ .

-  $N$  a set obtained by random independent  $m$  draws

$$m \geq \max \left( \frac{4}{3} \log \frac{4}{\delta}, \frac{8\delta}{\varepsilon} \log \frac{16}{\varepsilon} \right)$$

Then  $N$  is a  $C$ -net with prob.  $1 - \delta$

Remark:- Both of the thms. hold for spaces with shattering dim  $\delta$

$$\Theta \left( \frac{1}{\varepsilon} \log \frac{1}{\delta} + \frac{\delta}{\varepsilon} \log \frac{8}{\varepsilon} \right)$$

## Range Searching

Points in a d dimensional space  $\mathbb{R}^d$ .

We have a database.

Given a hyper-rectangle we want to report the points that are inside

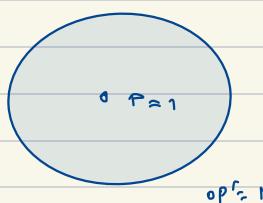
Allow, 1% error

# 6-sample-thm. says that there is a subset of constant size (which depends on  $\epsilon$ )  
use this to perform an estimation

Rects. have  $\mathcal{O}(1)$ -VC dimension

Sample with prob.  $1 - \delta$ .

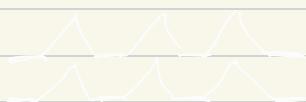
## Learning Concepts



Query oracle  
Assume, we know a function that returns  
"1" if inside, & "0" otherwise

→ There is a distribution  $D$  defined over the space.

We pick pts. from  $D$



## Growth-function

VC-dim = d

$$G_d(n) = \sum_{i=0}^d \binom{n}{i} \leq \sum_{i=0}^d \frac{n^i}{i!} \leq n^d$$

Lauer's Lemma:  $S: (X, R)$

Suppose VC-dimension of S  $\leq d$

$$\begin{aligned} |R| &\leq G_d(n) \leq \sum_{i=0}^d \binom{n}{i} \\ &\downarrow \\ \pi_f(n) \end{aligned}$$

Proof:- By induction on n, d. n=d=1 holds  
 Assume that it holds for n-1 & d-1  
 as well as n-1 & d-1.

We prove for n & d

$$\text{Define } f := \sum_{i=0}^d \binom{n}{i} = h(n, d)$$

Our induction hypothesis is for f with VC-dimension  $\leq d$

$$\pi_f(n) \leq d$$

$$\binom{n}{i} = \binom{n-1}{i} + \binom{n-1}{i-1}$$

Easy to verify  $h(n, d) = h(n-1, d) + h(n-1, d-1)$

recurrence,

Now let's fix a class  $F$

$\text{VC-Dim}(F) = d$  and a set  $X_1 = \{x_1, \dots, x_m\} \subseteq X$

$$F_1 = F|_{X_1}$$

$$F_2 = F|_{X_2}$$

$$F_3 = \{f|_{X_2} \mid f \in F \text{ and } f' \in F \text{ s.t. } \forall x \in X_2, f'(x) = f(x) \wedge f'(x_1) = -f(x_1)\}$$

$$\text{VC-Dim}(F') \leq \text{VC-dim}(F) \leq d$$

$$|F_1| = \frac{|F_2|}{\leq d} + \frac{|F_3|}{\leq d-1}$$

$$\frac{|F_1|}{|F|} \leq h(n-1, d)$$

$$|F_3| \leq h(n-1, d-1)$$

$$\Rightarrow |F| \leq h(n-1, d) + h(n-1, d-1)$$

$$\leq h(n, d)$$

Ex :- Let  $F$  be s.t.

$$\text{VC Dim}(F) \leq d \text{ for } n \geq d$$

$$\pi_F(n) \leq \left(\frac{n}{d}\right)^d$$

## Set Cover / Hitting Set (Piercing)

$U$  = Universe of elements

$X$  = Set of subsets

$S = (U, X)$  - set system

Choose a subset  $X' \subseteq X$  which is a cover - NP Hard

Greedy approx:

1. Sort the set based on cardinalities
2. Choose the set with max cardinality.

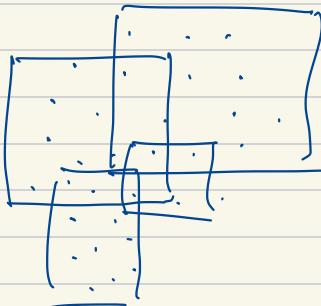
$\Rightarrow$  log-factor

$\hookrightarrow$  if a lower-bound shows that we can't get better than log factor.

wish for "nice" set families, can we bound Greedy (log-factor)

$S = (X, R)$   
↓  
Set of Elements       $\curvearrowright$  Set of Ranges

Goal: Choose a subset  $R' \subseteq R$  that covers  $X$  - minimum



Class of objects - has bounded VC-Dimension

↳ This implies that there is an epsilon-net

↳ Epsilon-net is a sample that "hits" all the heavy sets ( $\geq \varepsilon n$ )

A set  $N \subseteq \mathcal{X}$  is an  $\varepsilon$ -net for a finite subset  $n$ , if for any range

$$x \in R, m(x) = \frac{m(x)}{|x|} \geq \varepsilon \text{ then } N \text{ contains}$$

### Construction of $\varepsilon$ -net

1 - Choose a "random" sample. If it's large enough it's  $\varepsilon$ -net.

"Small hitting set" (which is also a net)

### $\varepsilon$ -net them

We can get a subset  $N$ , by  $m$  ind. draws for a finite subset  $n$

$$m \geq \max \left\{ \frac{1}{\varepsilon} \log \frac{4}{\delta}, \frac{8\delta}{\varepsilon} \log \frac{10}{\varepsilon} \right\}$$

with prob.  $> 1-\delta$

# Suppose the shattering dimension is  $\dim(d)$

$$\text{Sample size } \geq O\left(\frac{d}{\varepsilon} \log \frac{d}{\varepsilon}\right)$$

## Weighted Net

Suppose the elements are weighted ( $w: \mathcal{X} \rightarrow \mathbb{R}^+$ )

$\mathcal{L}$ -subset

$$w(\mathcal{L}) = \sum_{j \in \mathcal{L}} w(j)$$

Goal: Hit all the  $x_i$  with  $wt. \geq \varepsilon x_i w_i$ .



Algorithm  $\mathcal{L} = (\mathcal{X}, \mathcal{R})$  dual -  $\mathcal{L}^* = (\mathcal{X}^*, \mathcal{R}^*)$

1. Repeatedly select an  $\varepsilon$ -net (for some  $\varepsilon$ ) of  $\mathcal{L}^*$   $\rightarrow$  shattering dim
2. Size of the net =  $O\left(\frac{\delta^*}{\varepsilon} \log \frac{\delta^*}{\varepsilon}\right)$

3. Verify if it is a net. If not, discard.

↳ if it is a net - check if it is a set cover

If yes  $\rightarrow$  Done

If no  $\rightarrow$

- Let  $R_p = \{r \in \mathcal{R} \mid p \in r\}$  all ranges that contain

$p$

- Double the weight (Multiplicative weight update) of the elements in  $R_p$ .

Obs:- Every time we double, we are increasing not more than  $(1+\varepsilon)$  multiplicative factor

$$w_i \leq (1+\varepsilon)^i w_0$$

i many iterations.

Minimum weight of elements  $k = |\text{opt}|$  in opt?

$$\begin{aligned} k \times 2^{ik} &\leq w_i = (1+\varepsilon)^i w_0 = (1+\varepsilon)^i x m \\ &\leq e^{i\varepsilon} x m \end{aligned}$$

$$i = k \times g$$

$$k \times 2^g \leq e^{i\varepsilon} x m$$

$$\log k + g \leq \log m + \varepsilon i$$

$$g(1-\varepsilon k) \leq \frac{\log m}{k}$$

$$\text{Suppose we take } \varepsilon = \frac{1}{2k} \Rightarrow g \leq O\left(\log\left(\frac{m}{k}\right)\right)$$

The problem is that we don't know what  $\varepsilon$  is?

↪ Binary search  $\rightarrow$  this is almost a  $\log$  overhead to for whole process.

$\rightarrow$  Size of our core

$$O(s \times k \log(s \times k))$$



$k$  opt cover

Q How to choose  $\xi$ ?

$\xi$  is independent of  $k$ .  
Suppose we choose  $\xi = \frac{1}{4k}$ .

Guess we fix value of  $\xi$

$$O\left(k \lg \frac{m}{k}\right)$$

---

### Parametrized algorithms

 Hard Problems & Exact Solutions

Idea: We aim for exact algorithm

But we want to isolate  $\exp^{\Theta}$  term (parameters)

$\Rightarrow$  obtain very fast  $\exp^{\Theta}$ , when parameters small.

(Hope: parameters are small in practice)

Parameters — non-negative integer  $k(x)$  comes with probability

$\nu_p$ )

denote by  $k$

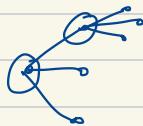
Not necessarily efficiently computable.

## Parametrized Problem

↳ Problem + Parameter ( $k$ )  
(wrt  $k$ )

Goal :- polynomial complexity on  $n$ .  
exponential complexity on  $k$ .

Example :-



I/p:  $G_i = (V_i, E_i)$ ,  $k \in \mathbb{N}$

O/p: Does there exist

a  $k$ -size VC (output a set  $S(\leq k)$  s.t.  $\forall e \in E, \exists v \in S \rightarrow v \in e$ )

## Brute force soln

- Try all  $\binom{n}{k} + \binom{n}{k-1} + \binom{n}{k-2} + \dots + \binom{n}{0}$ ,
- All sets of  $k$ -vertices
- Testing invalid VC takes  $O(E)$
- Total -  $O(V^k E)$  - poly. for fixed  $k$  - this is slow

## Branching (Bounded Search tree technique)

↳ pick an arbitrary edge  $e \in E$   
↳ We know either  $u \in S$  or  $v \in S$   
or  $\{u, v\} \subseteq S$

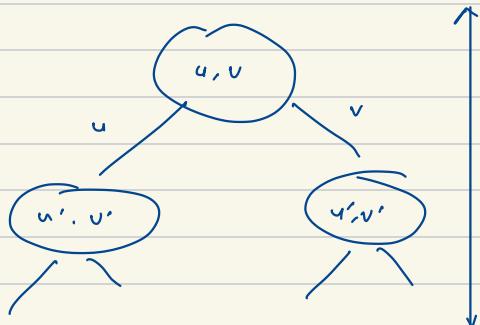
Guess - Try both

① Add  $u$  to  $S$

(delete  $u \in N(u)$  from  $S$ ). recursive  $K' = K-1$

② Same for  $v$

- Return OR of the both



Bottom-up  
 $\downarrow k$

$\Rightarrow O(2^k \times v)$

as long as  $k \leq \log v$

this is good.

### Fixed-parameter Tractable (FPT)

If  $\exists$  an edge with running time  $\leq f(k) \times n^{o(1)}$

$\hookrightarrow$  polynomial

$$f: \mathbb{N} \rightarrow \mathbb{N}$$

Q Why is  $f(k) \times n^{o(1)}$  and not  $f(k) + n^{o(1)}$ ?

Thm:-  $f(k) \times n^{o(1)} \Leftrightarrow f(k) + n^{o(1)}$

Proof:-  $\Rightarrow$  if  $n \leq f(k)$

$$f(k) \times n^c < f(k)^{ct+1}$$

: if  $f(k) \leq n$  then

$$f(k) \times n^c \leq n^{ct+1}$$

$$\text{So, } f(k) \times n^c \leq \max\{f(k)^{ct}, n^{ct}\}$$

$$\leq f(k)^{ct} + \frac{n^{ct}}{n^c}$$

$\Leftarrow$  Trivial assumption  
 $f(k) \approx n^{c'} \geq 1$

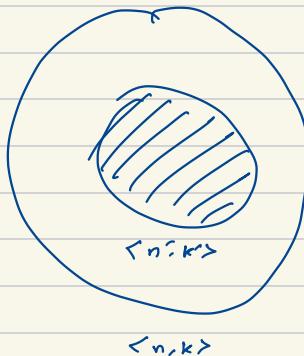
## Kernelization

$\rightarrow$  Simplifying self-reduction  
 (poly-time reduction)

$\rightarrow i/p \rightarrow \langle n, k \rangle$   $\Leftarrow \Rightarrow$   
 converts it into  $\langle n', k' \rangle$

How small :-  $|n'| \leq f(k)$

Equivalent :-  $\text{Ans.}(\langle n, k \rangle) = \text{Ans.}(\langle n', k' \rangle)$



Thm.

FPT  $\Leftrightarrow$  Kernelization

( $\Leftarrow$ ) Kernelize  $\Rightarrow n' \leq f(k)$  run any fpt  $g(w) \Rightarrow n^{O(1)} + g(f(k))$   
 time

( $\Rightarrow$ ) A runs in  $f(k) \times n^c$  if  $n \leq f(k)$  - kernelized

If  $f(k) \leq n$

- run A  $\Rightarrow f(k) \times n^c \leq n^{ct}$

O/P  $O(1)$  size Yes/No

## Sunflower Lemma (Erdos - Rado Conj)

- Classical results from 1960
- Apply in Kernelization

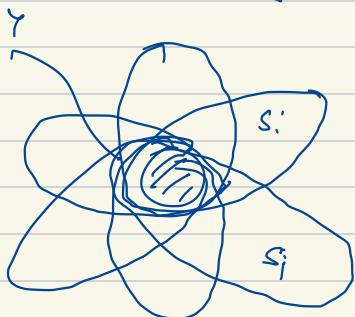
Sunflower is a collection of sets  $S_1, S_2, \dots, S_k$

s.t.  $S_i \cap S_j = Y \Leftrightarrow i \neq j \rightarrow$  Petals  $\forall i, S_i \setminus Y$

-  $K$  petals

- A core  $Y$

(None of them can be empty)



Lemma -  $F$  - family of sets (no duplication) over a universe  $U$ .

s.t. each set has cardinality exactly  $d$ .

- if  $|F| > d! (k!)^d$ , then

$F$  contains  $K$  petals

→ Poly-time to compute this

↳  $|F|, |U|, k$

Proof:- By induction on  $d$ , for  $d=1$ , Singaltons

Suppose  $d \geq 2$ , Let  $G = \{S_1, S_2, \dots, S_d\} \subseteq F$  be

inclusion-wise maximal family of pairwise disjoint sets  
in  $F$ .

$\rightarrow$  If  $l \geq k$ , then  $G$  is a sunflower already with at least  $k$  petals

$\rightarrow$  Assume  $l < k$ ,

$$S = \bigcup_{i=1}^l G_i$$

$$\text{Then } |S| \leq d \times (k-1)$$

Since  $G$  is maximal, every set  $A \in F$ , intersects at least one set from  $G$   $A \cap S \neq \emptyset$

Therefore, there is an element  $u \in U$ , which is contained in at least

$$\frac{|F|}{|S|} \geq \frac{d! (k-1)^d}{d(k-1)} = (d-1)! (k-1)^{d-1}$$

Set from  $F$

$\rightarrow$  Take all sets of  $F$  containing this element  $u$

Construct  $F'$  of sets with cardinality  $(d-1)$  by removing  $u$

$$|F'| \geq d! (k-1)^d$$

By i.h.

$F'$  contains sunflower  $\{s'_1, \dots, s'_{k'}\}$  with  $k'$  petals

$$\{\{s'_1 \cup u\}, \{s'_2 \cup u\}, \dots, \{s'_{k'} \cup u\}\}$$

## Poly time Alg.

→ Greedily select maximal sets  
if size is atleast  $R$ -done  
Else, find  $u$  and recurse.

## Erdos - Rado - 1960 (London Math Soc.)

Each set in  $F$  has cardinality  $d$ . If  $|F| \geq d!(R-1)^d$ , then  
there is a sunflower

$(\log k)^d$  (Annals. of Math 21, Alweiss et al.)

Terry Tao's blog - Sunflower Lemma via Shannon's  
entropy

## $d$ -Hitting Set (Application of Sunflower Lemma)

I/p → Family of Sets  $\mathcal{A}$  over  $u$   
each set has cardinality atmost  $d$ ,  
a non-negative integer  $k$

O/p → whether there is a subset  $H \subseteq u$  of size atmost  $k$   
st.  $H$  contains atleast one element from each set.

Proof :-

If  $A$  contains a sunflower,

Say  $S = \{s_1, s_2, \dots, s_{kn}\}$  of #  $kn$

then every hitting set  $H$  of  $A$  of cardinality at most  $k$  intersects its core (some  $y$ )

### Reduction Rule

$(U, d, k) \rightarrow \text{Return } (U', A, k)$

$$A' = (A \setminus S) \cup \{y\}$$

$$\Delta \quad U' = \bigcup_{x \in A} x$$

### Kernel formation algo.

If # sets  $\geq d! \times k^d$ . find a sunflower

- Apply recursively

→ not having it is  
not a problem since  
it's already  
~ kernel

Kernel size  $\leq d! \times k^d$ .

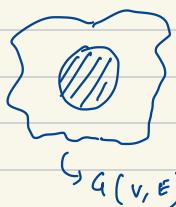
Fixed-parameter algo??

## Randomized Methods in Parametrized Algorithms

Color-coding [Alon, Yuster, Zwick J. ACM, 1995]

↳ Detect small subgraph in a large input graph

Framework:



I/p = n vertex graph  $G$

O/p = n vertex subgraph  $H$  of  $G$

(of some particular pattern)  
(clique, ind-set)

Brute-force  $\Rightarrow$  check every subgraph

↳  $O(n^k)$

Goal -  $2^{O(k)} * n^{O(1)}$

$\Rightarrow$  This is achievable for trees

# Trees - Treewidth (graphs which are not trees but tree-like)

Q Can we always get something?

- Most likely, no

$k$ -clique

Framework

Randomly color  $G$ , in such a way that a "pattern" emerges w.h.p.

## $k$ -simple path

(Remark - it is possible to de-randomize using splitters or pseudo-random generators)

$\text{I/P} - A$  graph  $G = (V, E)$ , tree integer  $k$   
 $O/P - k$ -simple path

Goal - Color vertices in  $V$ , uniformly at random from  $\{1, 2, \dots, k\}$  and find a path, if it exists

Hope - this process gives a "colorful" path.

Q There is a set  $S$  of  $k$  vertices. What is prob. that all vertices get different color?

$$\frac{k!}{k^k} \geq \frac{1}{e^k} \quad (\text{stirling approx})$$

## Algorithm

Repeat  $e^k \times t$  times

- Pick random  $f : V(G) \rightarrow \{1, 2, \dots, k\}$
- Look for a colorful path

Success - If the algorithm gets a path, then  $G$  has a path

Failure - If  $\exists$  a  $k$ -path but the alg. doesn't find it

$$1 - \left(1 - \frac{1}{e^k}\right)^{e^k \times t} \geq 1 - \frac{t}{e^t}$$

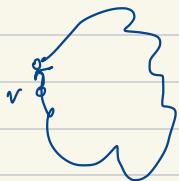
$\xrightarrow{\text{t}} \xrightarrow{\text{to success in one round}} \xrightarrow{\text{failure in one round}}$

# This algorithm doesn't give us the path

Remark: No randomization in colorful path finding

Homework:  $k^{O(k)}$  time alg. for finding colorful path.  
 $(8^{k+1}/O(k))$

→ Dynamic Programming based on simple observations



$T[S, v] = \begin{cases} \text{True} & \text{if } \exists \text{ a path of } |S| \text{ vertices ending in } v \\ & \text{using all elements of } S \\ \text{False} & \text{otherwise} \end{cases}$

$$\rightarrow \bigcup_{n \in N(v)} T[S \setminus f(v), n]$$

↳ # of table entries

$$= 2^k \times n^2$$

↳  $O(n)$  time for filling up one entry

$$\# \text{ of runs} = O^k \times t$$

$$\text{Total time} \rightarrow O((2c)^k \times n^2)$$

## Derandomization

↳ De-randomize

Q How can we make the random coloring deterministic?

- Let  $F = f_1, f_2, \dots, f_k$  be a family of  $f^2$

with  $f_i : V(G) \rightarrow \{1, 2, \dots, k\}$

→  $F$  is a universal hash family if for every set

$S \subseteq V(G)$  of size  $k$  there is  $f \in F$  s.t.  $f$  makes  $S$  colorful.

→ Construct a  $k$ -perfect hash family  $F$ .

For each  $f \in F$ , look for a colorful  $k$ -path

Total time -  $O(t + |F| \times 2^k \times n^2)$

↳ Naor, Schulman, Srinivasan,

[FOCS-1995]

↳ Construct  $k$ -perfect hash families  
of size  $e^{kt+O(t)} / \gamma n$