



CS 602

Approximation Algorithm

Design algorithm that strictly runs in polynomial time ($n^{O(1)}$)
 Output is allowed to be a "provable" factor away from
 the optimal solution.

Maximization Problems

Ind Let
 Variable $\alpha > 1$ set of vertices such that no two of them are connected

α -approx if we output a solution that is $(\frac{1}{\alpha})$ to the optimal solution

Minimization problems

Hamiltonian Cycle
 Cycle that visits every vertex of G exactly once and returns back α -opt if we output a solution that is at most α of the optimal solution

Polynomial-time approximation solution

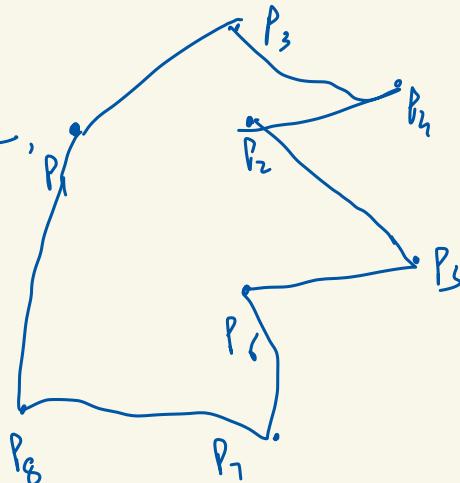
algorithm (of with some parameter $\epsilon > 0$) for any input, output a solⁿ within a factor $(1 + \epsilon)$ of the optimal solution.
 that runs in $n^{f(\frac{1}{\epsilon})}$ for some comparable f^n .
 (Running time is polynomial in $n, \text{some } \epsilon$)

Traveling Salesman Problem (TSP)

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 every city exactly once.

Decision Version

Given a length L ,
is it possible
to find a solution
of length L .



Graph:

A set of vertices, edges, weights $G = (V, E, w)$
 Visit all the vertices without repetition (minimize the sum
 of edge weights)

↓
Hamiltonian cycle problem

- * Hamiltonian cycle is NP-complete (Richard Karp)
- * No constant factor abs! is possible

Symmetric

Some edge weights
 on both
 directions

Asymmetric



Metric Space :

- $d(u, v) \geq 0$
- $d(x, y) = d(y, x)$

- Triangle inequality
 $d(x, y) + d(y, z) \geq d(x, z)$

x . y
 . z

$\text{cost}(S) = \text{sum of the weights of the edges (union)}$

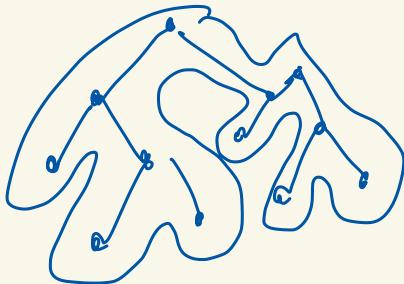
$S = \text{set of edges}$

same as finding the \min of Hamiltonian path
path cycle

Base Structure

- Min Spanning Tree [kruskal]

Due to triangular inequality, we can remove duplicates and cost is reduced



Do the DFS traversal and delete duplicates

Analysis

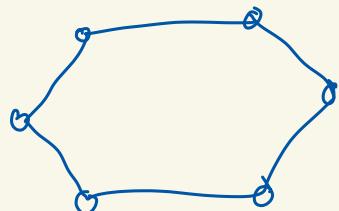
Every edge of the MST is traversed twice

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

valid cycle

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

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



Q Can we do better?

$$I/O = G = (V, E) \xrightarrow{\quad} OPT_G$$

(i) Take a subset

Induced subgroup

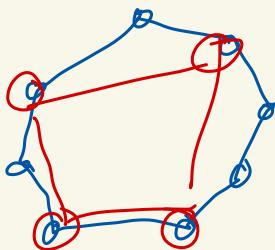
$$S \subseteq V$$

$$G_i(S)$$

$$\xrightarrow{\quad} OPT_S$$

$$OPT_S \leq OPT_G$$

Triangle
Inequality



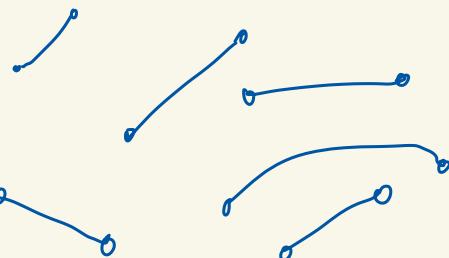
Property-2

Perfect Matching (can be computed in polynomial time)

Min cost AM

Perfect Matching with
smallest cost

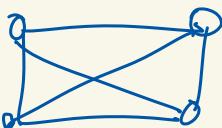
(Polynomial Time
 $\sim O(n^2)$)



Eulerian Tour (Circuit)

vertices can be repeated

each node has
even degree than
this is possible



$$\sum d(v_i) = 2 \times [\epsilon]$$

↑
every edge connected twice

Q How many odd degree vertices we have?
even

We will add another edge for perfect matching

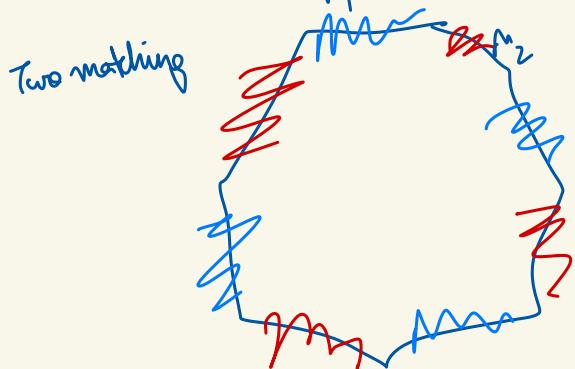
Compute Eulerian circuit

→ Vertex repetition is allowed, but we can delete the duplicates.

Analysis

$$\begin{aligned} \text{Cost } (C) &= \underbrace{\text{Cost } (\text{MST})}_{\leq \text{Cost } (\text{optimal})} + \text{Cost } (\text{Matching}) \\ \text{Cost } (C') &\leq \frac{1}{2} \text{Cost } (\text{optimal}) \end{aligned}$$

$\text{Cost } (C') \leq \frac{3}{2} \text{Cost } (\text{optimal})$



$$\begin{aligned} \text{Cost } (M_1) &\leq \text{Cost } (\text{optimal}) \\ \text{Cost } (M_2) &\leq \text{Cost } (\text{optimal}) \\ \text{Cost } (M') &\leq \frac{1}{2} \text{Cost } (\text{optimal}) \\ \text{Matching we choose because it is minimum} \end{aligned}$$

Metric TSP

- 2- α_{PR} [MST doubling]
- 1.5- α_{PR} (Christofides algorithm, 1976)

$$\downarrow \\ 1.5 - \varepsilon$$

$$\varepsilon = 10^{-30} \quad (2021)$$

No α_{PR} is possible if the distances are arbitrary.

$$G = (V, E, W)$$

- Determine if there is a hamiltonian cycle & some length t .

G' - complete graph

TSP in G' of wt n .

Does there exist in PTAS $(1+\varepsilon)-\alpha_{\text{PR}}$ for metric TSP
No $n^{o(1)}$ time

Theorem: There can't be a PTAS $(220/219) - \alpha_{\text{PR}}$, unless $P = NP$.

Restrict the metric

$\overline{\mathcal{I}}$

Euclidean metric

A set of points in \mathbb{R}^2 , with euclidean distances

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

Find the shortest route that covers all the points.

The Traveling Salesman Problem

$$\text{Cities} = \{1, 2, \dots, n\}$$

$C(n \times n)$ matrix \rightarrow Cost of traveling between pairs of cities

\downarrow
symmetric

$\underbrace{\quad}_{\text{If we view this as undirected complete graph}}$
then the problem is $\underbrace{\text{Hamiltonian cycle}}_{\text{problem.}}$

Approximation algorithms for the
TSP can be used to solve the Hamiltonian cycle problem

NP-complete

$$G_1 = (V, E)$$

$$C_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \\ n+2 & \text{otherwise} \end{cases}$$

If Hamiltonian cycle then
tour = n

$$\text{otherwise } \geq (n+2) + (n-1) = 2n+1$$

$\underbrace{\quad}_{\text{Input to TSP-algo}}$

we can detect
hamiltonian cycle $\leftarrow \begin{cases} 2\text{-apx can increase the cost to } 2n \\ \text{for hamiltonian cycle} \end{cases}$

\downarrow
Contradiction! (because HC is NP-complete)

Assumption: Restrict attention to metric space (metric TSP)

Algo(1): A spanning tree of a connected graph $G_1 = (V, E)$ is a minimal subset of edges $F \subseteq E$ such that each pair of nodes in G is connected by a path using edges only in F .

minimum spanning tree: Total edge cost minimized.

* Cost (optimal tour of TSP) \geq cost (MST)

Take this tour and remove one edge

(We will get a spanning tree whose cost $>$ cost (MST))

Algo(1) = nearest addition algorithm \rightarrow 2-apx algo

\downarrow

$F = \{(i_2, j_2), \dots, (i_n, j_n)\}$ \rightarrow edges obtained

$OPT > \sum_{l=2}^n c_{i_l j_l}$

\downarrow
Minimum spanning tree

Cost of the first
two nodes (i_2, j_2)

\downarrow
 $2c_{i_2 j_2}$ (traversed
two times)

j is inserted between (i, k)

marked

$$c_{ij} + \underbrace{c_{jk} - c_{ik}}_{\leq c_{ij}} \leq 2c_j$$

$$\text{cost (nearest-addition algo)} \leq 2 \sum_{l=2}^n c_{i_l j_l} \leq 2(OPT)$$

* Eulerian graph \rightarrow traversal of edges (each edge exactly once)

A graph is eulerian iff it is connected and each node has even degree

Algo(II) \rightarrow Double Tree Algorithm

MST compute \rightarrow replace each edge by two copies of itself

↓
resulting graph is Eulerian and has cost $\leq 2(\text{OPT})$

Eulerian Traversal \rightarrow sequence of edges (but vertices might repeat)

i_0, i_1, \dots, i_k remove all but the first occurrence of each city in this sequence.

↳ Tour of each city once

two consecutive cities (i_1, i_m)

we have removed i_{l+1}, \dots, i_{m-1}

By triangle inequality, cost is decreased,

↳ In total cost is at most the total cost of all the edges in the Eulerian graph
 $\leq 2(\text{OPT})$

double-tree = 2-apx algo.

Christofides Algorithm : MST Comput

↳ $O = \text{set of odd-degree vertices}$

For a tree, sum of degrees = $2 \times |E| = \text{even}$

↳ number of odd degree vertices = $|O| = \text{odd}$

$|O| = 2k$

→ perfect matching $(i_1, i_2), \dots, (i_{2k-1}, i_{2k})$

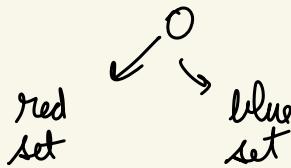
perfect-matching of minimum cost = $O(n^2)$

* Christofides = $\frac{3}{2}$ -apx algo

MST has cost $\leq \text{OPT}$

tour on nodes of $\underbrace{\mathcal{O}}$ has cost $\leq \text{OPT}$
↓
subset of original graph

Consider the shortest tour on the node set \mathcal{O} . Colour edges
red and blue



$$\text{cost(red)} + \text{cost(blue)} \leq \text{OPT}$$

$$\min(\text{cost(red)}, \text{cost(blue)}) \leq \frac{\text{OPT}}{2}$$

Perfect matching cost \leq Perfect matching + MST $\leq \frac{3}{2} \text{OPT}$

* For any constant $\alpha < \frac{220}{219}$ no α -apx for the metric TSP.

Euclidean TSP

Given n points in \mathbb{R}^2 with Euclidean distances i.e.

$$d(x_i, y) = \|x - y\|_2 \quad \text{shortest tour that visits all points?}$$

Euclidean TSP = NP-hard (do not know NP) length might be irrational

ϵ -nice instance

- (1) Every point has integral coordinates in the interval $[0, O(\frac{m}{\epsilon})]^2$
- (2) Any two different points have distances at least 4.

Consider the smallest bounding box around the points of the input instance. L = longer side of the box

$$L = \left\lceil \frac{8n}{\epsilon} \right\rceil$$

* OPT_I = length of the optimum tour in I . We can transform I into an ϵ -nice instance I' such that $OPT_{I'} \leq (1+\epsilon)OPT_I$

Proof: $I \rightarrow$ smallest bounding box \rightarrow length longer = L

optimal tour $\geq 2L$ \leftarrow $\begin{cases} \text{two points opposite in} \\ \text{the box have distance} \\ \geq L \end{cases}$

Now to obtain $I' \rightarrow$ draw a fine grid with spacing $\frac{\epsilon L}{2n}$ and map each point to the closest grid point.

$$L = \left\lceil \frac{8n}{\epsilon} \right\rceil = O\left(\frac{n}{\epsilon}\right) \quad \frac{\epsilon L}{2n} > \frac{\epsilon}{2n} \times \frac{8n}{\epsilon} = 4$$

$\Rightarrow I' = \epsilon$ -nice (integer coordinates and $d_{ij} > 4$)

By mapping points of I' to the points in I , we moved each point at most by $\frac{\varepsilon L}{2n}$

→ Edge changed by $\frac{\varepsilon L}{n}$ n edges $\rightarrow \varepsilon L$
 $\leq \text{opt}$

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

VC dimension

Range space $S = (X, R)$

elements of $X \rightarrow$ points
elements of $R \rightarrow$ ranges

↓
ground set
(finite or infinite)

family of subsets of X
(finite or infinite)

$x =$ finite subset of X

measure of a range $\bar{m}(x) = \frac{|x \cap X|}{|x|}$

subset N (might be a multi-set)
of x

estimate of the measure $\bar{m}(x)$
is $\bar{s}(x) = \frac{|x \cap N|}{|N|}$

$Y \subseteq X$ $R_{|Y} = \{x \cap Y \mid x \in R\}$

projections of R on Y . The range space
S projected to Y is $S_{|Y} = (Y, R_{|Y})$

If $R_{|Y}$ contains all subsets of Y ($\text{if } Y = \text{finite}, |R_{|Y}| = 2^{|Y|}$)

then Y is shattered by R

VC dimension ($\dim_{VC}(S)$) maximum cardinality of a
shattered subset of X .

Interval $\rightarrow VC = 2$

Disks $\rightarrow VC = 3$

Convex sets $\rightarrow VC = \infty$

Complement : range space $S = (X, R)$ $\dim_{VC}(S) = \bar{S}$

$$\bar{S} = (X, \bar{R})$$

$$\bar{R} = \{X \setminus \sigma \mid \sigma \in R\}$$

If S shatters B , then for any $Z \subseteq B$, $(B \setminus Z) \in R_{|B}$

$$Z = B \setminus (B \setminus Z) \in \bar{R}_{|B}$$

$\Rightarrow \bar{R}_{|B}$ contains all the subsets of B .

$\Rightarrow \bar{S}$ shatters $B \Rightarrow \dim_{VC}(\bar{S}) = \dim_{VC}(S)$

* Let $P = \{p_1, \dots, p_{d+2}\}$ be a set of $d+2$ points in \mathbb{R}^d . There are real numbers $\beta_1, \dots, \beta_{d+2}$ not all of them zero such that $\sum_i \beta_i p_i = 0$ and $\sum_i \beta_i = 0$

Proof: $q_i = (p_i, 1)$ $q_1, \dots, q_{d+2} \in \mathbb{R}^{d+1}$ are linearly dependent.

There are coefficients $\beta_1, \dots, \beta_{d+2}$ not all of them zero such that $\sum_{i=1}^{d+2} \beta_i q_i = 0$ considering only the first d -coordinates $\sum_{i=1}^{d+2} \beta_i p_i = 0$ $(d+1)^{th}$ coordinate $\sum_{i=1}^{d+2} \beta_i = 0$

Rado's Thm: $P = \{p_1, \dots, p_{d+2}\} \subset \mathbb{R}^d$ Then, there exist two disjoint subsets C and D of P , such that $CH(C) \cap CH(D) = \emptyset$

$$C \cup D = P$$

Proof: By previous thm, $\sum_i \beta_i p_i = 0$ and $\sum_i \beta_i = 0$

$$\mu = \sum_{i=1}^k \beta_i = - \sum_{i=k+1}^{d+2} \beta_i$$

$$\sum_{i=1}^k \beta_i p_i = - \sum_{i=k+1}^n \beta_i p_i$$

$v = \sum_{i=1}^n (\beta_i / \mu) p_i$ is a point in Convex Hull($p_1 \dots p_n$)

$$v = \sum_{i=k+1}^{d+2} -(\beta_i / \mu) p_i \in \text{CH}(p_{k+1}, \dots, p_{d+2})$$

v = intersection of the two convex hulls

* $P \subseteq \mathbb{R}^d$ = finite set s = point in $\text{CH}(P)$ h^+ = halfspace containing s . Then there exists a point of P contained inside h^+ .

Proof: $h^+ = \{t \in \mathbb{R}^d \mid \langle t, v \rangle \leq c\}$

$$\sum_i \alpha_i = 1 \quad \text{and} \quad \sum_i \alpha_i p_i = s$$

$$\langle s, v \rangle \leq c \Rightarrow \left\langle \sum_{i=1}^m \alpha_i p_i, v \right\rangle \leq c \Rightarrow \beta = \sum_{i=1}^m \alpha_i \langle p_i, v \rangle \leq c$$

$\beta_i = \langle p_i, v \rangle$ β is a weighted average of $\beta_1 \dots \beta_m$

\Rightarrow there must be a β_i which is no larger than the average $\Rightarrow \beta_i \leq c \Rightarrow \langle p_i, v \rangle \leq c \Rightarrow p_i \in h^+$.

* Growth Function = $G_S(n) = \sum_{i=0}^n \binom{n}{i} \leq \sum_{i=0}^n \frac{n^i}{i!} \leq n^S$

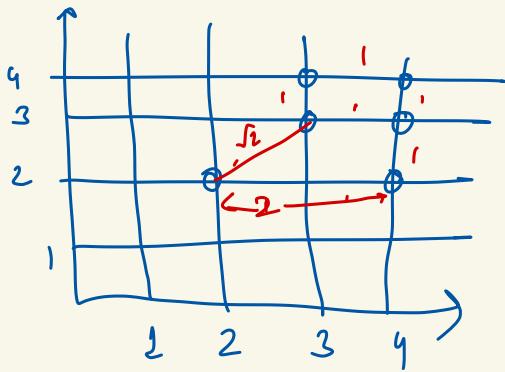
* Sauer's Lemma: If (X, R) is a range space of VC dimension S with $|X| = n$ then $|R| \leq G_S(n)$

Proof: holds for $n=0$ and $S=0$

$$R_x = \{\sigma \setminus \{x\} \mid \sigma \cup \{x\} \in R \text{ and } \sigma \setminus \{x\} \in R\}$$

$$R \setminus x = \{\sigma \setminus \{x\} \mid \sigma \in R\}$$

$$|R| = |R_x| + |R \setminus x| \leq G_{S-1}(n-1) + G_S(n-1) = G_S(n)$$



Euclidean TSP is NP-hard
but not known to be
NP

Sum of square roots (SRS)

Given a set of positive integers
 $\{a_1, \dots, a_k\}$ decide
 $\sum_{i=1}^k a_i \leq t$

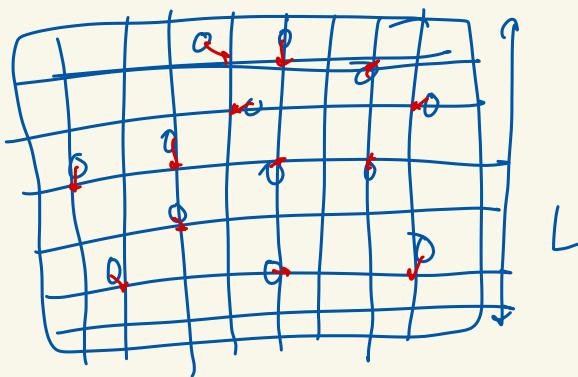
$$\{a_1, \dots, a_k\} \quad \{b_1, \dots, b_k\}$$

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

NP-hard
not in NP

PTAS - polynomial time approximation scheme.
for euclidean TSP

- Rounding the instance
- Partitioning (Exploit the structure of the instance by breaking it into more instances)
- Dynamic programming



Map each point
to the closest
grid point
(To make the
coordinates
rational)

Given - E

ϵ - "nice" instance

Defⁿ - An instance of Euclidean TSP is ϵ -nice if

1. Every point has integral co-ordinates in the interval $[0, O(\frac{n}{\epsilon})^2]$
2. Any two diff points have dist at least 4.

- Take a small bounding box (axis-parallel)

longer side - L.

s.t. rooted not origin

$$\text{Scale } L = \sqrt{\frac{8n}{\epsilon}}$$

Lemma - I is slp & OPT_I is optimal tour.
I' is ϵ -nice instance $OPT_{I'}$ is optimal tour

$$OPT_I \leq (1 + \epsilon) OPT_{I'}$$

OPT is at least $2L$

- Draw a fine grid with spacing $\frac{\epsilon \times L}{2n}$
- Map every pt to its closest grid point (multiple pts could be mapped to the same grid pt).
- All pts have integer coordinates

$$L = \left\lceil \frac{8n}{\epsilon} \right\rceil \text{ or } O\left(\frac{n}{\epsilon}\right)$$

$$\text{Grid spacing } \frac{\epsilon L}{2n} \Rightarrow \frac{\epsilon L}{2n} \times \frac{8n}{\epsilon} = 4$$

→ Mapping each point in I has moved $\frac{\epsilon L}{2n}$

Every edge in the sol'^m changes by at most

$$\frac{\epsilon L}{2n} \text{ edges in OPT}$$

$$\text{Cost} = \epsilon \times L$$

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

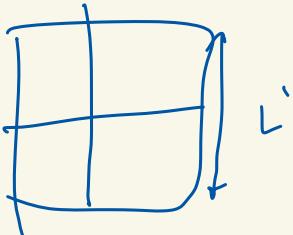
$$\text{OPT}_I + \epsilon \times \text{OPT}_I$$

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

$$L \leq \text{OPT}_I$$

Partition the space

- Extend the bonding box to square with new side length L'
 L' is the smallest power 2



* Recursively partition the box
 into four equal sized squares
 until the side length is 1
 ↴ (L' is power of 2)

- each pt is separated
- one pt in each "non-empty" square

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

Height of the quadtree

- $O(\log L')$
- $O(\lg(\frac{n}{\epsilon}))$

Apply dynamic programming to the quad tree

Solve for each square that are leaves

↓
 Bottom-up combine

Portals

Limits the # of iterations

of portals
 Accuracy improve
 Running Time } Tradeoff

Select $m = \text{power of 2}$

$$m \in \left[\frac{k}{\epsilon}, \frac{2k}{\epsilon} \right]$$

for each square \rightarrow put portals in corners
 put $(m-1)$ portals equally spaced

Portal-respecting tour (p-tour)

Defn: p-tour enters/exits through portals

$$-\text{ length of p-tour} \leq (1+\epsilon) \times \text{OPT}$$

Detours can add much more cost

Sol^m → Randomize

(i) Translate the grid by a random offset at most

$$\frac{1}{2} \text{ in each coordinate}$$

(ii) points remain grid points.

(iii) with high probability, the pts are nicely concentrated.

(iv) higher levels in the partition (quad tree)
have more portals → fine-grained tree

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

Theorem: (a,b) picked up uniformly at random $\left[0, \frac{1}{2}\right]$ with prob at least $\left(\frac{1}{2}\right)$

p-tour such that cost (p-tour) $\leq (1+4\epsilon) \times \text{OPT}$.

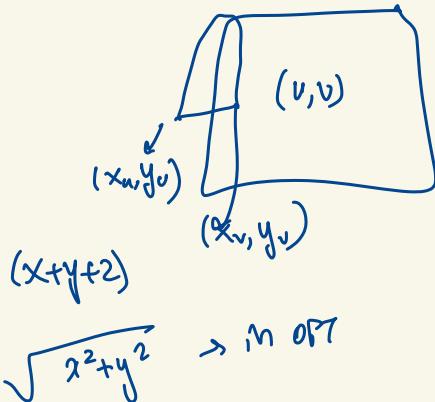
Extend non-ptour to a ptour

Proof: For each vertical/horizontal line l
 $+ (l) = \# \text{ of times intersects } l$

$$T_L = \sum_L t(L)$$

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

- e crosses $(x+1)$ vertical lines
- " " $(y+1)$ horizontal "
- total contribution



$$\begin{aligned} \sqrt{2(a^2+b^2)} &\geq a+b \\ \forall x,y \quad d(x,y) &\geq 0 \\ x+y+2 &\leq \sqrt{2(x^2+y^2)} + 2 \\ &\leq 2 \underbrace{\sqrt{x^2+y^2}}_{\text{OPT}} \end{aligned}$$

Bound - expected length of the detours

Detours might occur

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

i of the quad-tree

$$\frac{L'}{2^{i-m}}$$

if L is in level i

$$\leq \frac{L'}{2^{i-m}}$$

Q: what is the prob that after random shift L crosses L at level-i

- l could be mapped to $\frac{l'}{2}$ many times [translated by $(0, \frac{l'}{2})$]

- 2^{i-1} many lines of level i :

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

Expected length $\sum_{i=1}^k \frac{2^i}{l'} \times \frac{l'}{2^i m} \leq \epsilon$

By linearity of expectation $2\epsilon \times \text{OPT}$

Markov Inequality

Pr (total length

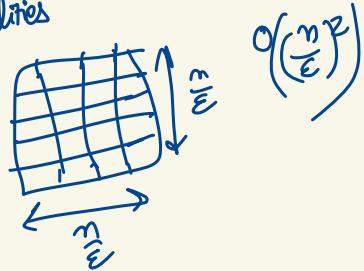
increase of detours $> 4\epsilon \text{ OPT}$)

$$\leq \frac{2\epsilon \text{ OPT}}{4\epsilon \text{ OPT}} = \frac{1}{2}$$

De-randomize

- fixed ϵ

- grid shifting by trying all possibilities



Final Step (DP)

Given (a, b) dissection, get p-tours

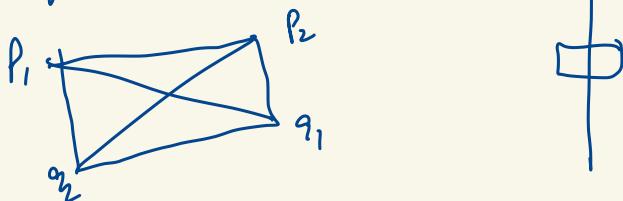
Introduce state

- Square

- any set of possible ways of entering/exiting the squares

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

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



- 4 portals
- use one portal $m = O\left(\frac{n}{\epsilon}\right) = O\left(\log \frac{L}{\epsilon}\right)$

$$\text{Catalan number} = \frac{1}{2n+1} \binom{2n}{n} = O(2^{2n}) = O(2^{kn})$$

Algorithm = try all parenthesis

- translate them into paths

- Discard anything that intersects

$$m = O\left(\log \frac{n}{\epsilon}\right) \quad \# \text{ of entry exits} = O(n^{1/\epsilon})$$

Computation of values

A $\left[(s_1, t_1) \dots (s_\ell, t_\ell) \right]$ - Compute the whole table.

Clustering

- Learning, searching, data mining
- Given data, find an interesting structure
- Represented as points in \mathbb{R}^d

General metric space (X, d) where X is a set
 $d : X \times X \rightarrow [0, \infty)$

is a metric it satisfies -

- (i) $x=y \rightarrow d_\mu(x, y) = 0$
- (ii) $\forall x, y \quad d_\mu(x, y) = d_\mu(y, x)$
- (iii) $\forall x, y, z \quad d_\mu(x, y) + d_\mu(y, z) \geq d_\mu(x, z)$

Assumption

$(x, y) \quad d_\mu(x, y) \quad$ in $O(1)$ time

Norm

↳ norm defines distances between pts
 $p, q \in \mathbb{R}^d \quad \|p-q\|_p = \left(\sum_{i=1}^d |p_i - q_i|^p \right)^{\frac{1}{p}}$ for $p \geq 1$

$p=2$: Euclidean norm

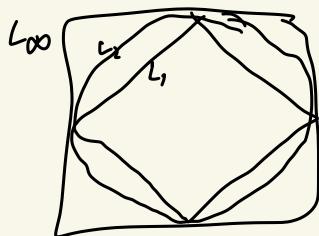
$p=1$: Manhattan distance (L_1 norm)

ℓ_∞ norm

$$\|p - q\|_\infty \leq \lim_{p \rightarrow \infty} \|p - q\|_p$$

max $|p_i - q_i|$

Triangle inequality holds for ℓ_∞ too, it's called Minkowski inequality



for any $p \in \mathbb{R}^d$

$$\|p\|_p \leq \|p\|_2 \quad \text{if } p \geq 0$$

Lemma — For any $p \in \mathbb{R}^d$

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

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

Const. a $f(x) = x^2 + (a-x)^2$ minimized if $x = \frac{a}{2}$

$$\text{Let } \alpha = \|p\|_1 = \sum_{i=1}^d |p_i|$$

By symmetry obs on $f(x) = \sum_{i=1}^d x_i^2$

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

Metric space (X, d)

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

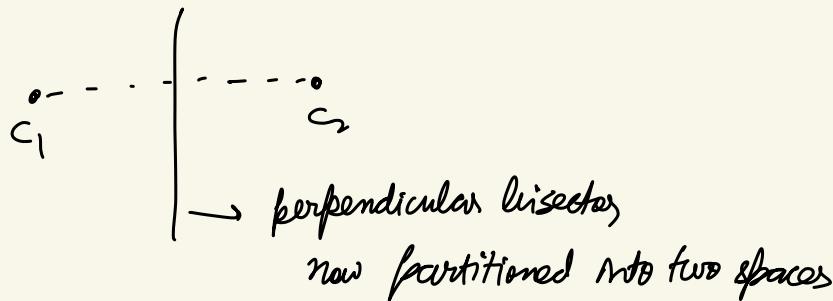
O/P : Find a clustering (set of centers) such that each pt is assigned to its nearest center

Set of clusters C

$$\text{Cluster}(C, \bar{c}) = \{p \in P \mid d_N(p, \bar{c}) = \underbrace{d_N(p, c)}_{\text{J}}$$

Voronoi Partition

minimum
across all the
pts



$$p_c = (d(p_1, c), d(p_2, c), \dots, d(p_n, c))$$

i^{th} coordinate $d(p_i, c)$ dist to p_i to its closest center

O/P = Find a set of k -centers $C \subseteq P$

such that the maximum distance of a point in P to its closest center is minimized

Defⁿ: Given a set of k -centers C ,
 $\|p_c\|_\infty = \max_{p \in P} d(p, C)$

Find C , s.t $\|p_c\|_\infty$ is minimized

$$\text{opt}_\infty(P, k) = \min_C \|p'_c\|_\infty \quad C \subseteq P \quad k = |C|$$

- C_{opt}
- NP-hard
- Hard to approximate beyond 1.86
- 2-approx in the Euclidean space in \mathbb{R}^2

Greedy Algo

- Start by picking an arbitrary pt. \bar{c}_1

$$C_1 = \{\bar{c}_1\}$$

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

- Take the pt with worst distance

$$(x_1 = \max_{p \in P} d_1(\{p\}))$$

say \bar{c}_2

$$C_2 = C_1 \cup \{\bar{c}_2\}$$

$$C_i = C_{i-1} \cup \{\bar{c}_i\}$$

$O(n)$
space

$O(nk)$
T
data from
previous
iterations

For each pt $p \in P$, a single variable $d[p]$ with its current dist to the closest pt.

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

Defⁿ: A ball of radius σ around a pt $p \in P$ is a set of pts in P with dist at most σ from p

$$b(p, \sigma) = \{q \in P \mid d_N(p, q) \leq \sigma\}$$

Remark: k -center is essentially covering P with k -balls of minimum radius.

Thm: Greedy Algo computes a set K of k -center such that k is 2 -apx $\|p_k\|_1 \leq 2 \|p_k\|_\infty$ takes $O(n \times k)$ time.

Proof: Running Time ✓

$$\text{Def}^n \quad r_k = \|p_k\|_\infty$$

Let \bar{c}_{k+1} is the point realising

$$r_k = \max_{p \in P} d(p, k)$$

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

By the defⁿ of r :

$$r_1 \geq r_2 \geq \dots \geq r_k$$

$$i < j < k+1$$

$$d_N(\bar{c}_i, \bar{c}_j) \geq d_N(\bar{c}_i, \bar{c}_{i-1})$$

$$r_{i-1} \geq r_k$$

— the dist. between any pair of pts. in C is at least τ_k

opt — covers P by using k balls

by triangle inequality any two points within such a ball are with a dist at most $2 \times \text{opt}$.

↓
None of the balls contain two points from
contradiction!

$$C \\ \subseteq P$$

Greedy permutation

Let this run till we exhaust all pts

$$\langle P \rangle = \langle \bar{c}_1, \bar{c}_2, \dots, \bar{c}_n \rangle$$

$$\downarrow \\ \langle \tau_1, \tau_2, \dots, \tau_n \rangle$$

Defⁿ : τ -packing : A set $S \subseteq P$ for P

(i) covering property : all the pts in P are within dist of atmost τ from S .

(ii) separation property : $\forall p, q \quad d_M(p, q) \geq \tau$

τ -packing gives compact representation

* Greedy permutation gives such a rep.

Thm: $\langle \overline{c_1}, \overline{c_2}, \dots, \overline{c_n} \rangle < \infty$
 for any i , we have $c_i = \langle \overline{c_1} \dots \overline{c_i} \rangle$
 is an ∞_i -packing of P

Proof: By contradiction

$$\infty_{k-1} = d(\overline{c_k}, \overline{c_{k-1}}) \forall k = 2, \dots, n$$

$$\text{for } j < k \leq i \leq n$$

$$d_\mu(\overline{c_j}, \overline{c_n}) = \infty_{k-1} \geq \infty_i$$

K-medians clustering

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 center is minimized.

Clustering price: $\|P_C\| = \sum_{p \in P} d(p, C)$

Objective f^* : $\text{opt}_p(p, k) = \min_{\substack{C \subseteq P \\ |C|=k}} \|P_C\|$

Optimal set of centres - C_{opt} .

Local search: move sol^n to sol^n in the space of candidate sol^n (the search space) by applying local changes

Continue until, end up on optimal or we exhaust the running time.

Notations:

$$\text{A set } U = \{P_c \mid C \in P^k\}$$

$$\text{opt}_{\infty}(P, k) = \min_{\substack{q \in U \\ \text{k-center}}} \|q\|_{\infty} \quad \left| \begin{array}{l} \text{opt}(P, k) = \min_{q \in U} \|q\|_1, \\ \text{k-median} \end{array} \right.$$

1.86 Apx X

2 Apx ✓ (Greedy)

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

$$\text{opt}_{\infty}(P, k) \leq \text{opt}_1(P, k) \leq n \times \text{opt}_{\infty}(P, k)$$

$$\begin{aligned} \text{Proof: } P &\in \mathbb{R}^m & \|P\|_{\infty} &= \max_{i=1}^n |P_i| \\ && \leq \sum_{i=1}^n \|P_i\|_1 &= \|P\|_1 \end{aligned}$$

$$\|P\|_1 \leq \sum_{i=1}^n |P_i| \leq \sum_{i=1}^n \max |P_i| \leq n \times \|P\|_{\infty}$$

C -set of centers $|C| = k$ realising $\text{opt}_1(P, k)$ i.e.

$$\text{opt}_c(P, k) = \|P_c\|_1$$

$$\begin{aligned} \text{opt}_{\infty}(P, k) &\leq \|P_c\|_{\infty} \\ &\leq \|P_c\|_1 = \text{opt}_c(P, k) \end{aligned}$$

Similarly, k realizing $\text{opt}_{\infty}(P, n)$

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

($2n$ -factor for median)

$2n$ -apx

use this as a first step for local search

L - is $2n$ -apx

Improve : parameter $0 < \delta < 1$
 $\forall i \in [n] L_{\text{curr}}$

Local search

- Set $L_{\text{curr}} \leftarrow L$
- At each iteration



We will check if the current "sol" L_{curr} can be improved

by replacing one of the centers

by one center from outside (non-centers)



Swap

$$K \leftarrow (L_{\text{curr}} \setminus \{\bar{c}\}) \cup \{\bar{c}\}$$

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

- continue swap as long as it satisfies the constraint
- return L_{curr}

Running time : An iteration takes $O(m \times k)$ swaps

$(n-k)$ candidates to be swapped in k -candidates
to be out)

implementing swap (naively $O(n^k)$)
overall $O(n^{2k})$

Since

$$\frac{1}{1-s} \geq (1+s)$$

$$O((n^k)^2 \log \frac{1}{1-s} \frac{\|P_k\|_1}{\epsilon_{\text{pt}_1}})$$

$$= O(n^k)^2 \cdot \log(1+s)^{2n} = O((n^k)^2 \log \frac{n}{s})$$

K-means

Set $P \subseteq X$, K , find K pts $C \subseteq P$ $|C|=k$

$$\|P_C\|_2^2 = \sum_{p \in P} (d_{\mu_p}(p, C))^2$$

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

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

$O(n)$ -factor for k -means as well

Thm: $0 < \varepsilon < 1$

$(25 + \varepsilon)$ -approx

VC-dim

- A range space $(X, R) = S$

X = ground set (finite / infinite)

R = family of subsets of X .

Consider finite subset of X as the estimating ground set.

Dfⁿ (Measure): fixed subset of X . For a range $\tau \in R$

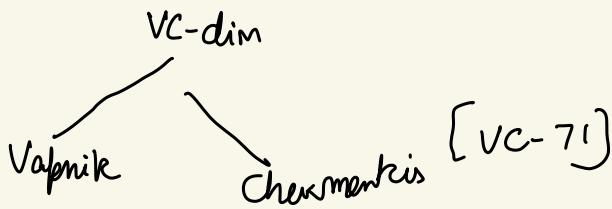
$$m(\tau) = \frac{|\tau \cap X|}{|X|}$$

For a subset N (multi-set) of X , the estimate of the measure of $m(\tau)$, for $\tau \in R$

$$\hat{s}(\tau) = \frac{|\tau \cap N|}{|N|}$$

Q2 How we get methods to generate N s.t.

$$\overline{S}(x) = \overline{m}(x) \quad \forall x \in \mathbb{R}$$



Dfm: $S = (X, R)$ For $Y \subseteq X$

$$R_{S,Y} = \{\tau \cap Y \mid \tau \in R\}$$

be the projection of R on Y

$\binom{|Y|}{2}$

If this is the cardinality then it is called
shattered by R

The orange
space S_Y
is projected
to $S_{Y'} = \{Y, R_{Y'}\}$

Complement

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

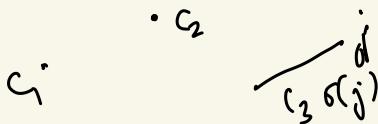
$$\overline{S} = (X, \overline{R}) \quad \text{where } \overline{R} = \{X \setminus \tau \mid \tau \in R\}$$

Q2 what is the VC-dim of \overline{S} ?

A subset $B \subseteq X$, is shattered in \overline{S} iff it is shattered in S
for any $Z \subseteq B$ $(B \setminus Z) \in \overline{R}_{IB} \Rightarrow Z = B \setminus (B \setminus Z) \in \overline{R}_{IB}$

Local search

- X be the set of arbitrary subset of k -centers
while true do:
(Swap) if i.e X and $i' \in F \setminus X$
 $\text{cost}(X - i + i') < \text{cost}(X)$
- (greedy solution of k -centers)



$\text{opt} = X^* - \text{optimal set of } X \leftarrow X - i + i'$
k-center otherwise break

Nearest

$$i^* \in X^*$$

$$i \notin X$$

for each center chosen in X ,
nearest center in X^* $\min_{i^*}(\|i\|_j)$

Inverse

Ties the centers in X^*

inverse is the
nearest map

Bijection

Claim: for any $j \in X$,
 $d_{\text{nearest}}(\sigma^*(i), j) \leq d_j + 2d_j^*$

Half-spaces

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

Claim: $P = \{P_1, \dots, P_{d+2}\}$ set of points in \mathbb{R}^d

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

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

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

pts are linearly dependent . and these are
 $a_1, a_2, \dots, a_{d+2} \in \mathbb{R}^{d+2}$ coefficients $\beta_1, \dots, \beta_{d+2}$
 s.t. $\sum_{i=1}^{d+2} \beta_i p_i = 0$

- Considering first d -coordinates of these pts implies

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

$$\text{Similarly } (d+1) \text{ coordinates } \sum_{i=1}^{d+2} \beta_i = 0$$

Radon's Thm: $P = \{P_1 \dots P_{d+2}\}$ \exists disjoint subsets
 $C \& D$ of P . $H(C) \cap H(D) = \emptyset$ then
 $C \cup D = P$

Shattering Dim:

Property : A range space (\mathcal{R}) with $VC\text{-dim}(S)$
 means # of ranges given polynomially on (n)
 (Generally this is \exp^n)

$$\underline{\text{Growth function}}: \quad G_{\delta}(n) = \sum_{i=0}^{\delta} \binom{n}{i} \leq \sum_{i=0}^{\delta} \frac{n^i}{i!} \leq n^{\delta} \quad \text{for } \delta > 1$$

$$\underline{\text{Sauer's Lemma}}: \quad S = (X, R) \\ \text{VC}(S) = \delta \quad |X| = n \quad |R| \leq G_{\delta}(n)$$

Proof: $n=0 \quad \delta=0 \quad \rightarrow \text{done!}$

$$x \in X$$

$$\text{contains}_x \{R_x\} = \left\{ \sigma \setminus \{x\} \mid \sigma \cup \{x\} \in R \text{ and } \sigma \supset \{x\} \in R \right\}$$

$$\text{does not contain}_x \{R \setminus x\} = \left\{ \sigma \setminus \{x\} \mid \sigma \in R \right\}$$

$$\text{Observation: } |R_x| + |R \setminus x| = |R|$$

Shatter function: $S = (X, R)$ shatter f^n
 $\pi_S(m)$ is the maximum # of sets that might be created by S , when restricted to the subsets of size m .

$$\pi_S(m) = \max_{\substack{B \subseteq X \\ |B|=m}} |R_{|B|}|$$

Shattering dim: The smallest d such that $\pi_S(m) = O(m^d)$ $\forall m$

Then $S = (X, R)$ has shattering dim d , then the VC-dim is bounded by $O(d \log d)$

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

$$2^s = |R_{|N|}| \leq \pi_S(m)$$

$$s \leq \log c + d \log s \quad (s \geq \max(2, \frac{2}{c}))$$

$$\Rightarrow \frac{s \cdot \log c}{\log s} \leq d$$

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

$$f(x) = \frac{x}{\log x} \rightarrow \text{non-increasing } x > c$$

$c > \sqrt{e}$ if $f(x) \geq e$ $\Leftrightarrow x > 1$
 $f(x) \leq x$ then $x \leq \log x$

ε -net and ε -sampling

$S = (X, R)$ x is a finite subset of X
 $0 \leq \varepsilon \leq 1$

Informally, ε -sampling captures R , upto some ε -error

a subset $c \subseteq x$ is an ε -sample for x if for any range $r \in R$

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

\downarrow measure \curvearrowright estimate

Thm: (ε -sample, VCT₁) — There is a free constant C s.t. if (X, R) is any range space with VC dim S .

$x \subseteq X$ finite subset of X and $\forall \varepsilon, \phi > 0$
 \exists a random subset $C \subseteq X$ of

(with probability $= \phi$) cardinality $S = \frac{C}{\varepsilon^2} \left(\delta \log \frac{\delta}{\varepsilon} + \log \frac{1}{\phi} \right)$

ε -net : A set $N \subseteq X$ is an ε -net for x if for any range space $r \in R$ if $\bar{m}(r) \geq \varepsilon$
 then r contains at least one pt. of N (i.e. $r \cap N \neq \emptyset$)

(Intuitively: hit all heavy subsets)

ϵ -net Thm (HWF7)

$S = (x, R)$ has $\text{VC dim}(S) \leq S$

x is a finite subset of X .

Suppose $0 \leq \epsilon \leq 1$ & $\phi < 1$

- N a set obtained by random independent draws.
- $m \geq \max\left(\frac{4}{\epsilon} \log \frac{4}{\phi}, \frac{8S}{\epsilon} \log \frac{16}{\epsilon}\right)$
- Then, N is a ϵ -net with prob $(1-\phi)$.

* Remark: Both of the thms hold for spaces with shattering dim S . ($O\left(\frac{1}{\epsilon} \log \frac{1}{\phi} + \frac{S}{\epsilon} \log \frac{S}{\epsilon}\right)$)

Range Searching $p \in \mathbb{R}^d$ we have a database
Given a hyper rectangle, we want to report the points that lie inside

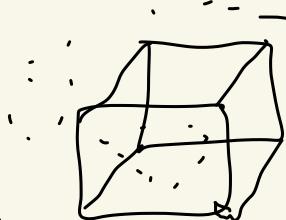
Allow 1% error,

ϵ -sample (Thm) says there is
a subset of const. size (which depends on ϵ)

Use this to perform an estimation.

Rectangle has bounded VC-dim

Random sample with probability $(1-\phi)$



Learning Concepts

Assume we know a f^* that returns 1 if inside,
0 otherwise

Query
Oracle

There is a distribution D defined over the space. We pick points from D .

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

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

Suppose $\text{VC dim}(S) \leq d$

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

$T_F(n)$

Proof: By induction on n, d
 $n=d=1$ holds

Assume that it holds for $n-1 \& d$
and as well as for $n-1 \& d-1$

We prove for $n \& d$
define $f^* : \sum_{i=0}^d \binom{n}{i} = h(n, d)$

Our induction hypothesis is for F with $\text{VC-dim} \leq d$

$$T_F(n) \leq d$$

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

$h(n, d) = h(n-1, d) + h(n-1, d-1)$

recurrence

Now let's fix a class F

$$VC\text{-dim}(F) = d \quad \text{and a set}$$

$$X_1 = \{x_1, \dots, x_m\} \subseteq X$$

$$f_1 = f_{1X}$$

$$f_2 = f_{2X}$$

$$F_3 = \{f_{1X} | f \in F \text{ & } f' \in F \text{ s.t. } \forall x \in X_2, f'(x) = f(x) \text{ & } f'(x_1) = -f(x_1)\}$$

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

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

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

Induction hypothesis

$$\begin{cases} |F_2| \leq h(n-1, d) \\ |F_3| \leq h(n-1, d-1) \end{cases} \rightarrow |F_1| \leq h(n-1, d) + h(n-1, d-1) \leq h(n, d)$$

Ex. Let F be s.t. $VC\text{-dim}(F) \leq d$ for $n \geq d$

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

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 approximation - (1) sort all the sets based on cardinality
(2) choose the set with max cardinality.

log factor

↳ \exists a lower bound shows that we can't get better than log factor.

Wish: for "nice" set families, can we beat greedy (log factor)?

$S = (X, R)$

↓
set of elements → set of ranges

Goal: choose a subset $R' \subseteq R$ that covers X .

class of objects \rightarrow has bounded VC-dim

ε -net: Sample that "wits" all the heavy sets ($> \varepsilon_n$)

A set $N \subseteq X$ is an ε -net for a finite subset x if
for any range $r \in R$, $m(r) = \frac{m(x)}{|x|} > \varepsilon$

then r contains at least one pt.

Construction of ε -net

- choose a random sample if it is large enough its ε -net
- small hitting set " (which is also a net)

ε -net thm

We can get a subset N by m ind. draws for a finite subset x .
(uniformly chosen)

$$N \geq \max \left\{ \frac{9}{\varepsilon} \log \frac{a}{\phi}, \frac{8\delta}{\varepsilon} \log \frac{16}{\varepsilon} \right\}$$

with prob $> 1 - \phi$.

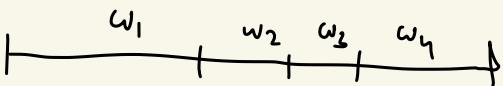
Suppose the shattering dim is d .

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

Weighted net: suppose the elements are wt,
($W: x \rightarrow \mathbb{R}^+$)

r -subset
 $w(r) = \sum_{j \in r} w(j)$

Goal : all the r 's with wt. $> \epsilon x_w$



Algorithm for set cover

(1) Repeatedly select an ϵ -net (for some ϵ)

$$S = (X, R) \text{ dual} - S^* = (X^*, R^*)$$

\downarrow
shattering dim S^*

$$\text{size of the net} = O\left(\frac{\xi^*}{\epsilon} \log \frac{\xi^*}{\epsilon}\right)$$

Verify if it is a net. If not - discard.

\downarrow
if it is a net check if it is a setcover
if yes - done

Let $R_p = \{\infty \in R \mid p \in \infty\}$ all ranges that contain p .

Double the weight of the elements in R_p

Observation: every time we double we are increasing not more than $(1+\epsilon)$ multiplicative function

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

$i \rightarrow$ iteration

Q1 What is the min wt. of elements ($K = \text{optt}$) in opt?

$$K \times 2^{\frac{i}{k}} \leq w_i = (1+\varepsilon)^i w_0 = (1+\varepsilon)^i x_m$$

$$\leq \varepsilon^{\frac{i}{k}} x_m$$

$$K \times 2^{\frac{i}{k}} \leq w_i$$

$$i = k \times g$$

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

$$\log(k + g) \leq \log m + \sum i$$

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

Suppose, we take $\varepsilon = \frac{1}{2k}$

$$\Rightarrow g \leq O\left(\log \frac{m}{k}\right)$$

$$\# \text{ of iteration } O\left(k \log \frac{m}{k}\right)$$

Size of our cover $O(S^*k \log S^*k)$ $k = \text{opt. cover}$

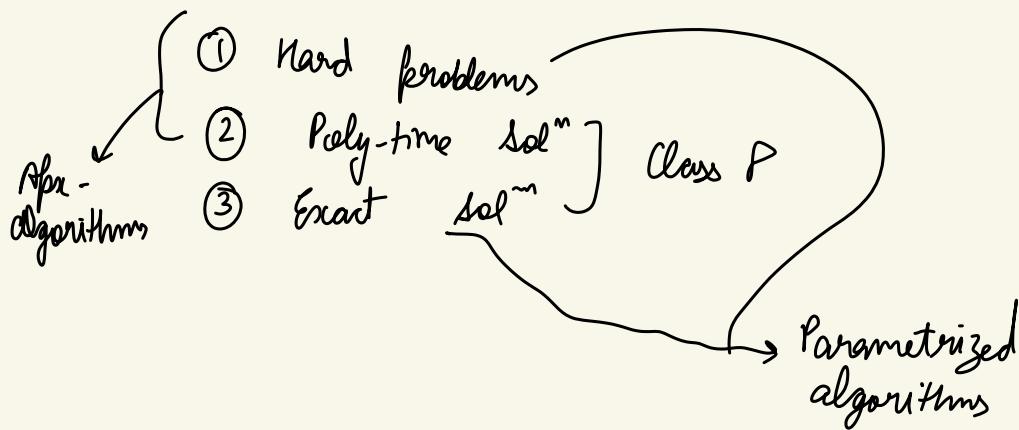
Q2 How to choose ε ?

ε is independent on k

Suppose $\varepsilon = \frac{1}{uk}$ instead of $\frac{1}{2k}$

Guess the value of $\frac{\varepsilon}{k_i}$

$O(k_i \log \frac{m}{k_i})$ iterations



Idea: Aim is to get exact algo

But we want to isolate \exp^n terms (parameters)

\Rightarrow obtain very fast sol'' when the parameter small.

(Note: parameters are small in practice).

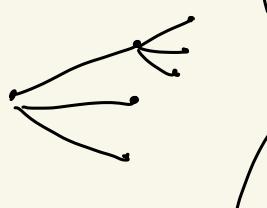
parameters - non-negative integer $k(x)$ (comes with prob i/p)
 - denote by k
 - Not necessarily efficiently computable.

Parametrized Problem

problem + parameter (k)
 (w.r.t. k)

Goal: poly. complexity on n
 Expⁿ complexity on k

Example



I/p : $G = (V, E)$ $k \in \mathbb{N}$

o/p : Does there exist a
 k -size vertex cover

↓
output a set ($\subseteq V$) s.t. $\forall e \in E$
 $\exists v \in S$

Brute force solution

(1) Try all $\binom{n}{k} + \binom{n}{k-1} + \dots + \binom{n}{0}$
 ↓
 All sets of k vertices

— Test valid VC takes $O(E)$ time

— Total = $O(V^k E)$ kely for fixed k .

slow for large n and reasonable k .

Branching (Bounded search tree technique)

→ Pick an arbitrary edge $e \in E$
 ↓
 (v, v')

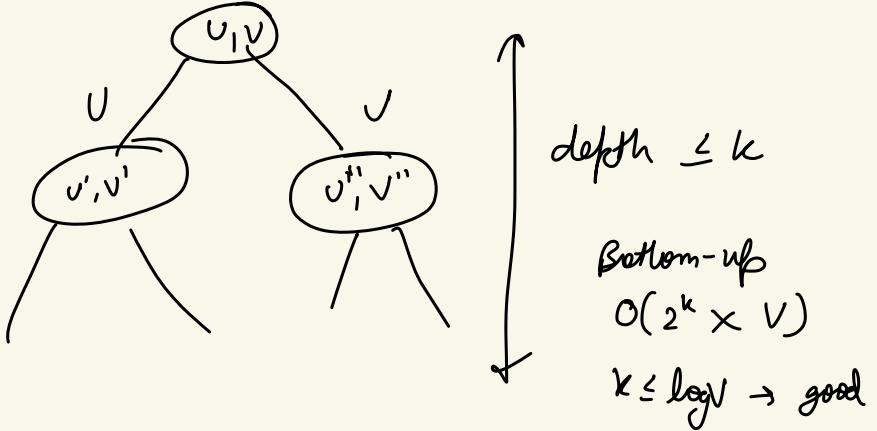
→ Know either $v \in S$ or $v' \in S$ or $\{v, v'\} \in S$

Guess — Try both

① Add v to S
 (delete v & $N(v)$ from S)
 recursive $k' = k - 1$

② same for v' .

— Return or of the outcomes



Fixed parameter-tractable (FPT)

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

$f: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ polynomial
 parameters

Q: why $f(k) \times n^{O(1)}$ and not $f(k) + n^{O(1)}$?

Thm: $f(k) \times n^c \Leftrightarrow f(k) + n^{c'}$

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

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

if $f(k) \leq n$

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

$$\text{So. } f(k) \times n^c \leq \max \left\{ f(k)^{c+1}, n^{c+1} \right\} \leq f(k)^{c+1} + \underbrace{n^{c+1}}$$

(\Leftarrow Trivial, assuming $f(k) & n^{c'} \geq 1$)

Kernelization

simplifying
self-reduction

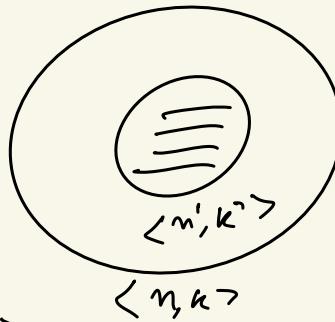
(poly-time reduction)

i/p - $\langle n, k \rangle$

converts it into $\langle n', k' \rangle$

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

Equivalent — Ans($\langle n, k \rangle$) = Ans($\langle n', k' \rangle$)



Thm:

FPT \Leftrightarrow kernelization

kernelization $\Rightarrow n' \leq f(k)$

run any finite $g(n')$

$\Rightarrow n^{O(1)} + g(f(k))$ time \rightarrow FPT

\Leftrightarrow

A runs in $f(k) > n^c$

if $n \leq f(k)$ in kernelized

if $f(k) \leq n$

run A $\rightarrow f(k) \times n^c \leq n^{c+1}$

O/p, O(1) size

k is known in advance

Sunflower Lemma (Erdős Rado Cons)

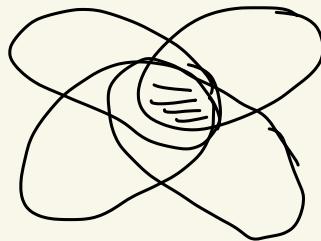
- Classical result from 1960

- Apply in kernelization

- k petals

- A core γ

(\Rightarrow None of them can be empty)



Collection of sets S_1, \dots, S_k

s.t. $S_i \cap S_j = \gamma$

$\forall i \neq j$

Petals $\forall i \in S \setminus \gamma$

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

s.t. each set has cardinality exactly d .

- If $|F| > d! (k-1)^d$, then F contains k petals

- Poly-time algorithm to compute this

w.r.t. $|F|, |U|, k$

Proof: for $d=1$, singletons

Suppose $d \geq 2$

Let $G = \{S_1, S_2, \dots, S_l\} \subseteq F$

- If $l \geq k$ then G is a sunflower

already with at least k petals

Assume $l < k$

be inclusion-wise
maximal family of
pairwise disjoint sets
in F

$$S = \bigcup_{i=1}^k S_i; \quad \text{Then } |S| \geq d \times (k-1)$$

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

There is an element $v \in V$ which is contained in at least

$$\frac{|F|}{|S|} \text{ sets.} \quad \frac{|F|}{|S|} > \frac{d! (k-1)^d}{d(k-1)} = \underbrace{(d-1)!}_{\text{sets for } F} (k-1)^{d-1}$$

- Take all sets of F containing this element v .

Construct F' of sets union cardinality $(d-1)$ by removing v .

$$|F'| > d!. (k-1)^d$$

By induction hypothesis

F' contains a sunflower $\{S'_1 \dots S'_n\}$ with k -petals

$$\{S'_1 \cup v\} \dots \{S'_n \cup v\}$$

Poly-Time Algorithm

(1) greedily select maximal sets. If size is at least k done. Else find v and return.

Erdos-Rado - FIGO

Each set in F has cardinality d

If $|F| > d!(k-1)^d$ then there is a sunflower.

$(\log k)^d \rightarrow$ recent bound.

d -Hitting set

(Application of Sunflower Lemma)

Input: Family of sets A over V . each set has cardinality at most d .

a non-negative integer k

Output: whether there is a subset $H \subseteq V$ of size at most k , such that H contains 1 element of each of sets of A .

Proof: If A contains a sunflower, say $S = \{S_1, \dots, S_{k+1}\}$ of $\#(k+1)$ then every hitting set H of A of cardinality at most k intersects its core Y .

Reduction rule : (V, A, k)

Return (V', A', k') $A' = (A \setminus S) \cup \{x\}$
and $V' = \bigcup_{x \in A'} X$

If # of sets are larger than $d! \times k^d$ find a sunflower
— Apply green consider kernel size $O(d! k^d)$

Kernelization

A data reduction rule for a parametrized problem \mathcal{Q} is a function $\phi : \Sigma^* \times \mathbb{N} \rightarrow \Sigma^* \times \mathbb{N}$ that maps an instance (I, k) of \mathcal{Q} to an equivalent instance (I', k') of \mathcal{Q} such that ϕ is computable in time polynomial in $|I|$ and k .

$$\text{size}_A(k) = \sup \left\{ |I'| + k' : (I', k') = A(I, k), I \in \Sigma^* \right\}$$

* A kernelization algorithm for a parametrized problem \mathcal{Q} is an algorithm A that, given an instance (I, k) of \mathcal{Q} works in polynomial time and returns an equivalent instance (I', k') of \mathcal{Q} . Moreover, we require that $\text{size}_A(k) \leq g(k)$ for some computable function $g : \mathbb{N} \rightarrow \mathbb{N}$.

* If a parametrized problem \mathcal{Q} is FPT then it admits a kernelization algorithm

Proof: $\mathcal{Q} = \text{FPT} \Rightarrow \exists A (I, k) \in \mathcal{Q}$ in time $f(k)|I|^c$
 (I, k) algo runs A on (I, k) for at most $|I|^{c+1}$ steps
 If it terminates with an answer, use that for yes/no.
 If A does not terminate within $|I|^{c+1}$ steps, then return (I, k) itself

$$f(k) \cdot |I|^c > |I|^{c+1} \Rightarrow |I| < f(k)$$

$$|I| + k \leq \underbrace{f(k) + k}_{\text{computable}} \quad (\text{kernel size})$$

Sunflower Lemma

A sunflower with k petals and a core γ is a collection of sets S_1, \dots, S_k such that $S_i \cap S_j = \gamma$ for all $i \neq j$; the sets $S_i \setminus \gamma$ are petals and we require none of them to be empty (γ can be empty).

* Let A be a family of sets (without duplicates) over a universe U , such that each set in A has cardinality exactly d . If $|A| > d!(k-1)^d$, then A contains a sunflower with k petals and such a sunflower can be computed in time polynomial in $|A|, |U|$ and k .

For $d=1$, family of singletons, statement holds
 $d \geq 2$ $A = \text{family of sets of cardinality at most } d \text{ over a universe } U \text{ such that } |A| > d!(k-1)^d$.
 $G = \{S_1, \dots, S_l\} \subseteq A$ be an inclusion-wise maximal family of pairwise disjoint sets in A .
If $l \geq k$ then G is a sunflower with at least k petals.

G is maximal, every set $A \in A$ intersects at least one set from G i.e. $A \cap S \neq \emptyset$.

$$S = \bigcup_{i=1}^l S_i \quad |S| \leq d(k-1)$$

There is an element $v \in V$ contained in at least

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

sets from A . We take all sets of A containing such an element v , and construct a family A' of sets of cardinality $d-1$ by removing from each set the element v . Because $|A'| \geq (d-1)! (k-1)^{d-1}$, by the induction hypothesis A' contains a sunflower $\{S'_1, \dots, S'_r\}$ with k -petals. Then $\{S'_1 \cup \{v\}, \dots, S'_k \cup \{v\}\}$ is a sunflower with k -petals.

d -hitting set

Given a family A of sets over a universe V , where each set in the family has cardinality at most d , and a positive integer k . The objective is to decide whether there is a subset $H \subseteq V$ of size at most k .

such that H contains at least one element from each set in A .

* d -Hitting sets admits a kernel with at most $d!k^d$ sets and at most $d!k^d \cdot d^2$ elements.

Let (U, A, k) be an instance of d -hitting set and assume that A contains a sunflower $S = \{S_1, \dots, S_{k+1}\}$ of cardinality $k+1$ with core y . Then return (U', A', k') where $A' = (A \setminus S) \cup \{y\}$ is obtained from A by deleting all sets $\{S_1, \dots, S_{k+1}\}$ and by adding a new set y and $U' = \bigcup_{X \in A'} X$.

Additional Notes

* Voronoi Partitions: set of centers C , every point of P assigned to nearest neighbour in C

$$\Pi(C, \bar{C}) = \{p \in P \mid d(p, \bar{C}) \leq d(p, c)\}$$

* Greedy clustering algorithm: arbitrary point \bar{c}_1 into C , for every point $p \in P$ compute $d_{\bar{c}_1}(p)$ from \bar{c}_1 . Pick point \bar{c}_2 with highest distance from \bar{c}_1 . Add this to the set of centers and denote this expanded set of centers as C_2 .

overall algorithm = $O(nk)$

→ This algorithm is 2-approx.

Proof: Case-1 Every cluster of C_{opt} contains exactly one point of k .

$$p \in P$$

$\bar{c} =$ center p belongs in C_{opt}

$\bar{k} =$ center of k that is in $\Pi(C_{opt}, \bar{c})$

$$d(p, \bar{c}) = d(p, C_{opt}) \leq r_{\infty}^{opt}(p, k)$$

$$d(\bar{k}, \bar{c}) = d(\bar{k}, C_{opt}) \leq r_{\infty}^{opt}$$

$$d(p, \bar{k}) \leq d(p, \bar{c}) + d(\bar{c}, \bar{k}) \leq 2r_{\infty}^{opt}$$

Case-2: Two centers \bar{k} and \bar{v} of k both in $\Pi(C_{opt}, \bar{c})$

σ was added later

$$r_{\infty}^k(p) \leq r_{\infty}^{C_{i-1}}(p) = d(\bar{v}, C_{i-1})$$

$$\leq d(\bar{v}, \bar{k})$$

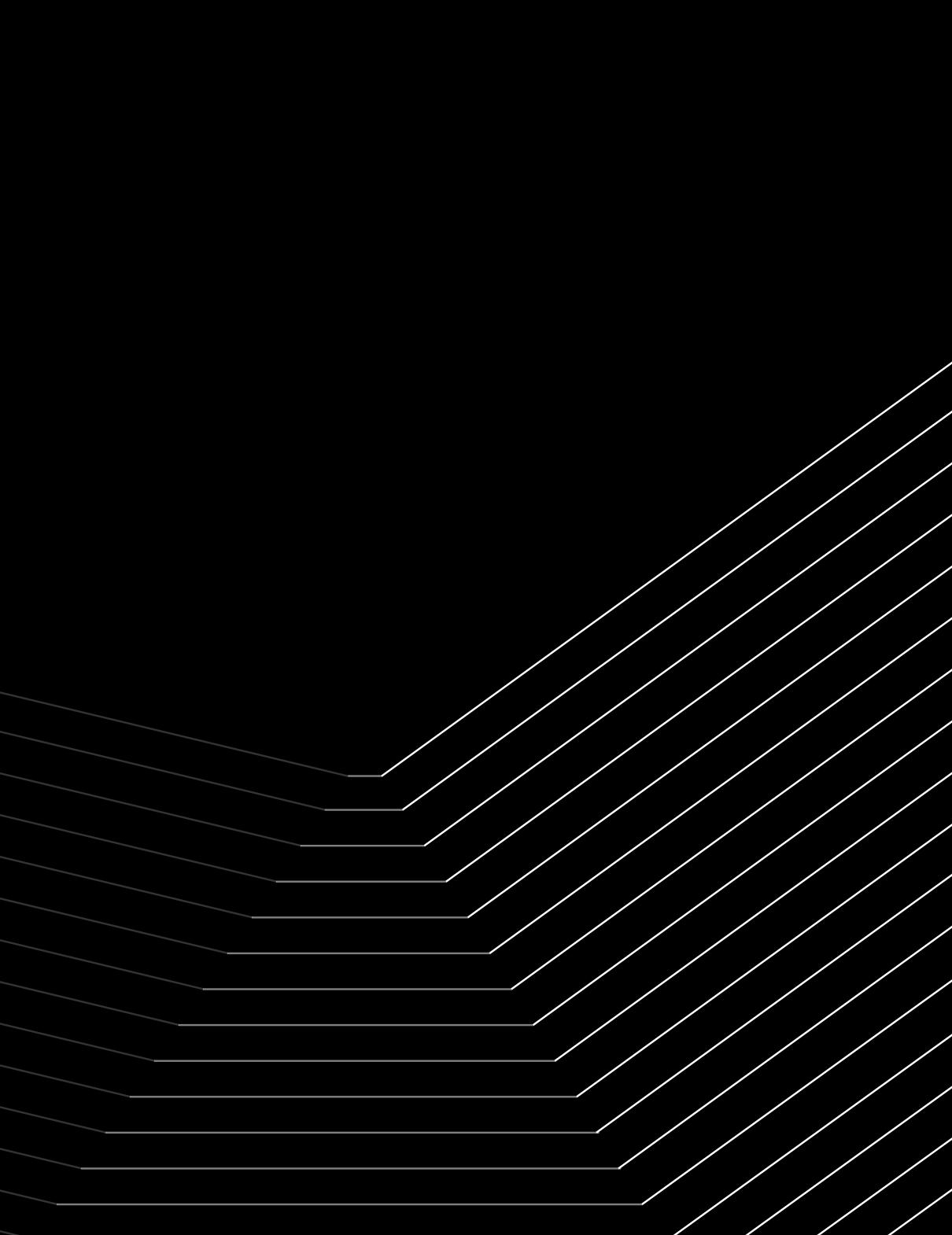
$$\leq d(\bar{v}, \bar{c}) + d(\bar{c}, \bar{k}) \leq 2r_{\infty}^{opt}$$

* A set $S \subseteq P$ is a σ -net for P if the following two properties hold :
 (i) Covering property = All the points of P are in distance at most σ from the points of S .
 (ii) Separation property = for any pair of points $p, q \in S$
 $d(p, q) \geq \sigma$.

* Let P be a set of n -points in a finite metric space, and let its greedy permutation be $\langle \bar{c}_1, \dots, \bar{c}_n \rangle$ with the associated sequence of radii $\langle \bar{\sigma}_1, \dots, \bar{\sigma}_n \rangle$ for any i , $C_i = \langle \bar{c}_1, \dots, \bar{c}_i \rangle$ is a σ_i -net of P .

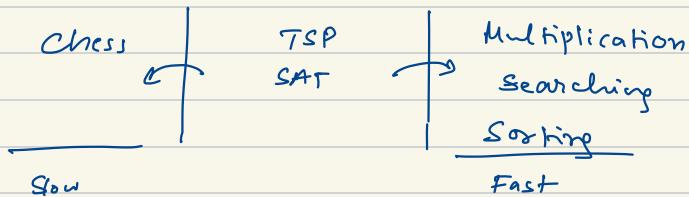
* $0 < p < 2 \quad \|x\|_p \geq \|x\|_2$

$$\|x\|_p \leq \sqrt{n} \|x\|_2 \quad \text{and} \quad \|x\|_2$$



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

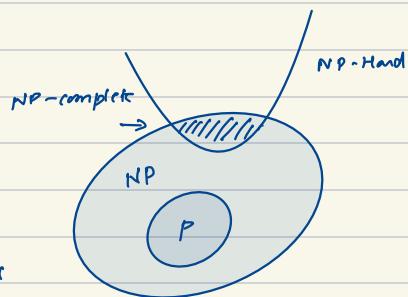
$NP = \left\{ \begin{array}{l} \text{Decision problems solvable in polynomial time} \\ \text{by non-deterministic TM} \end{array} \right\}$

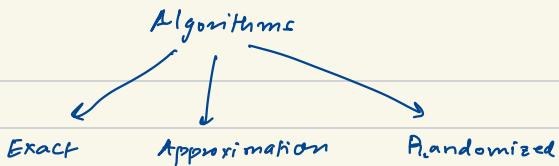
I/P \rightarrow [Verifier] \rightarrow Yes

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 $\alpha \geq 1$. α -approximation

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

Minimization Problems

Eg: Hamiltonian Cycle.

α -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 .

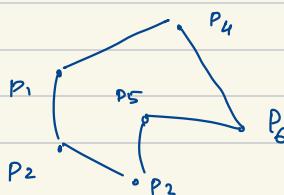
Approx. Algos - Vazirani, Shmoy's

(Running time is polynomial on n , some ε)

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

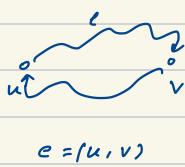
Minimizing the sum of edge weights

* Hamiltonian Cycle Problem

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

→ In fact, No constant factor approx is possible

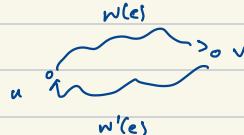
Symmetric



$$e = (u, v)$$

$$w(e).$$

Asymmetric



$$w(e)$$

Metric TSP

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



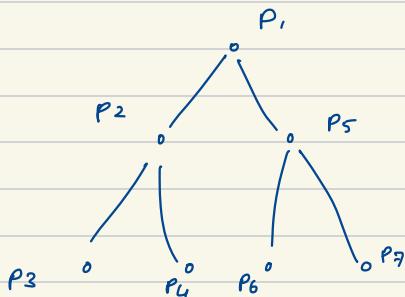
S - Set of Edges

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

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

Base structure:

- Min. Spanning Tree (Kruskal)



Repeat not allowed

P₁ P₂ P₃ P₂ P₄ P₁ P₅ P₆ P₅ P₇ P₅ P₁

↓
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

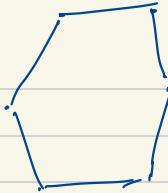
↳ Due to triangle inequality, Duplication is always possible

Analysis: Obviously, algorithm is polynomial time.

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

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

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



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?

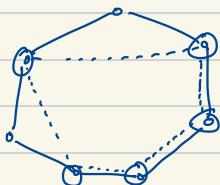
PROPERTY 1

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

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)

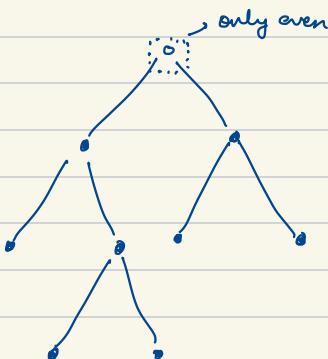


Minimum cost perfect matching is a perfect matching with smallest cost.
 (Poly.)

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:

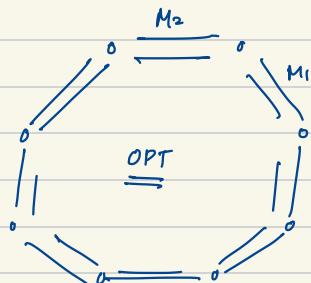
$$\text{Cost}(C) = \text{cost}(\text{MST}) + \text{cost}(\text{Matching})$$

↓ 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

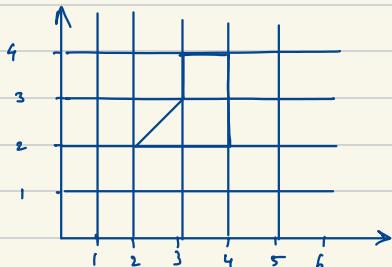
Thm: There can't be a PTAS $2^{20/21\alpha}$ -apx unless $P=NP$

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

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

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

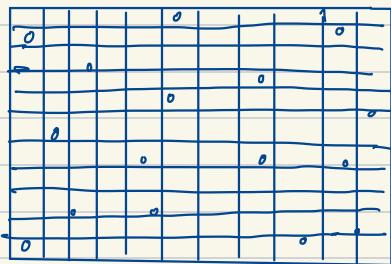
$$\sum_{i=1}^k a_i \leq \sum_{i=1}^k b_i$$

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"

ϵ - "nice" instance

Def² An instance of Euclidean TSP is ϵ -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 ϵ -nice instance $\Rightarrow \text{OPT}_{I'}$ is optimal tour

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

We know that OPT is at least ϵL ,

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

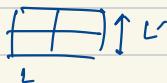
$$\leq \int_L^U L \leq \text{OPT}_I$$

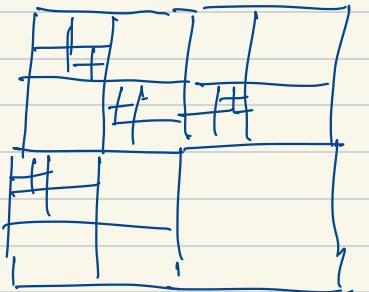
$$\text{OPT}_I + \epsilon \times \text{OPT}_I$$

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

② Partition the Space

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



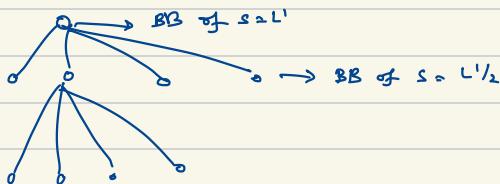


Quadtrees

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

At the end of this partition, we can guarantee that there will be atmost 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 $n^{O(\ell(\frac{1}{\epsilon}))}$

- ϵ -niceness

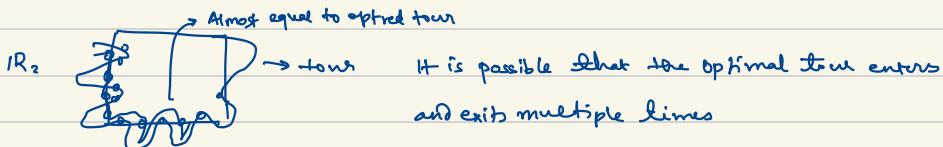
- quadtree partitioning

- DP

ϵ -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,

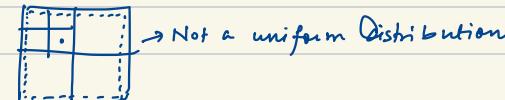
- Put portals in corners
- Put $(m-1)$ portals equally spaced
- Few many portals

Portal Respecting Tours

Def: p-tour enters & exits through portals

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

- Detours can add much more cost.



Sol¹: 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.

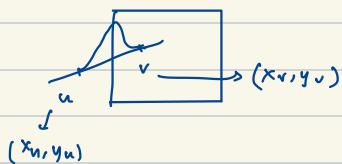
* For each vertical/horizontal line ℓ

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

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

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

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

By linearity of expectation

$$2\epsilon \text{OPT}$$

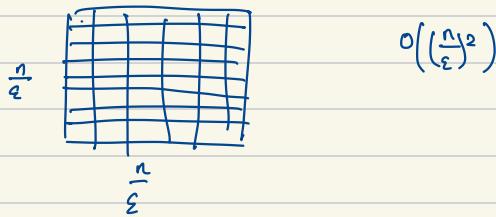
Markov Inequality

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

De randomization

- fixed ϵ

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

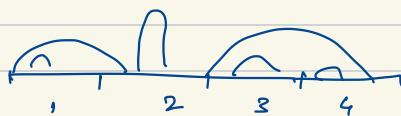
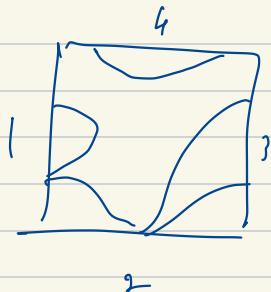
$$1 + 4 + 4^2 + \dots + L^2 = O\left(\frac{n^2}{\epsilon^2}\right)$$

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
 $p, q \in \mathbb{R}^d$

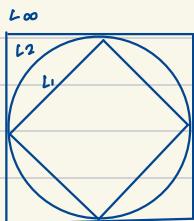
$$\|p - q\|_t = \left(\sum_{i=1}^d |p_i - q_i|^t \right)^{\frac{1}{t}} \quad \text{for } t \geq 1$$

$t=2$: 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_∞ 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. x

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

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

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

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

→ 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  $\pi_i$ ,

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

$$i < j < K+1$$

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

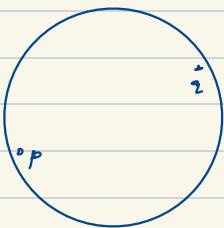
- the distance between any pair of points in  $C_i$  is at least  $r_k$

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  $r_k$  from  $s$ .

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

\*  $r$ -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  ${}^n C_k$  possible solutions - so an exhaustive local search should lead to the optimal  $\text{sol}^n$ .

### Notations

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

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

$k$ -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

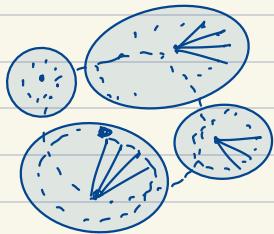
Improve - parameter  $0 < \delta < 1$

$w \in [n]$   
↳ curr

## Local Search -

→ Set :  $L_{curr} \leftarrow L$

→ At each iteration



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

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

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

→ 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  $0 < \varepsilon < 1$

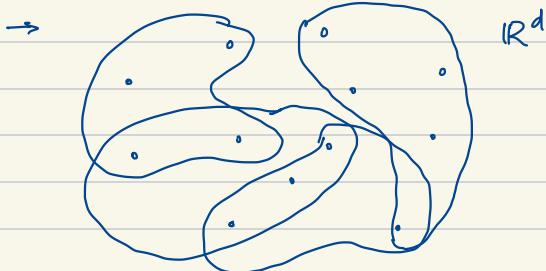
$$(25 + \varepsilon) \text{ factor appd } O(n^2 k^3 \frac{\log n}{\varepsilon})$$

## 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.

$\underline{X \setminus N}$

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

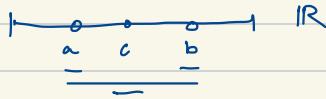
$$\hat{s}(x) = \hat{m}(x), \forall x \in \mathbb{R}$$

VC-Dimension  
↳ chernomskis  
Vapnik (1971)

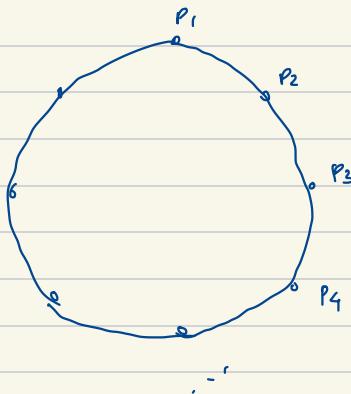
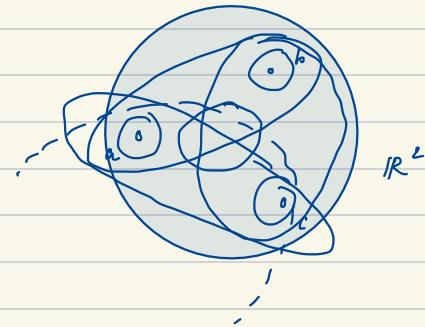
Def<sup>n</sup>  $S = (X, R)$  for  $Y \subseteq X$

$R_Y = \{s \cap Y \mid s \in R\}$  ↪ 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

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

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

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 :

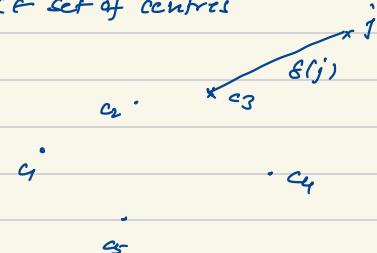
swap : if  $i \in X$  and  $i' \in F \setminus X$   
 $\text{cost}(x - i + i') < \text{cost}(x)$  :  
 $x \leftarrow x - i + i'$

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

$\min_{i' \in X^*} d(i, i')$

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)

Growth f. ..  $G_S(n) = \sum_{i=0}^{\delta} \binom{n}{i} \leq \sum_{i=0}^{\delta} \frac{n^i}{i!} \leq n^\delta$  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

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

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

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

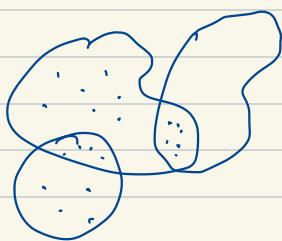
$$\frac{\delta}{2 \log \delta} \leq d \Rightarrow \frac{\delta}{\log \delta} \leq O(1 \times d)$$

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

- 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)

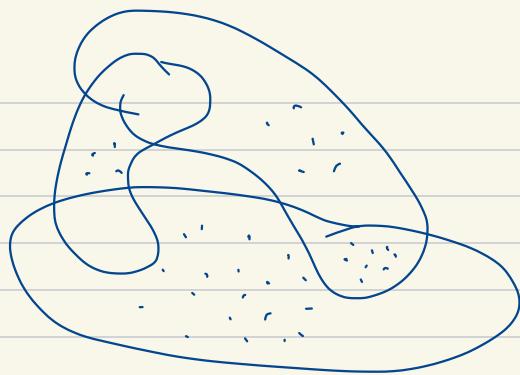
Theorem ( $\varepsilon$ -sample, VC Th): There is a tree constant  $c$ , s.t. if  $(X, R)$  is any range space with  $rc$ -dimension  $\delta$ .

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

a random subset  $C \subseteq X$  of cardinality

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

with probability  $1 - \phi$



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

$$|r \cap N| \geq \varepsilon |r|$$

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  $0 \leq \varepsilon \leq 1$  &  $\varphi < 1$ .

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

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

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

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

$$\Theta \left( \frac{L}{\varepsilon} \log \frac{L}{\varphi} + \frac{\varepsilon}{\varepsilon} \log \frac{8}{\varepsilon} \right)$$

## 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

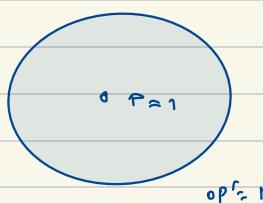
Allow, 1% error

# 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.  $1 - \delta$ .

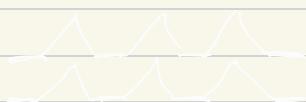
## Learning Concepts



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.

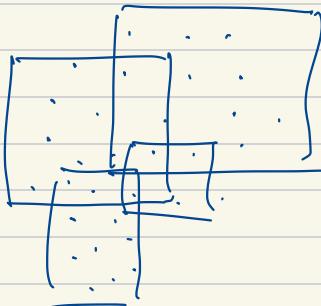
$\Rightarrow$  log-factor

$\hookrightarrow$  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  $k$

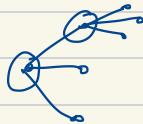
Not necessarily efficiently computable.

## Parametrized Problem

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

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

Example :-



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

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  $\{u, v\} \in S$

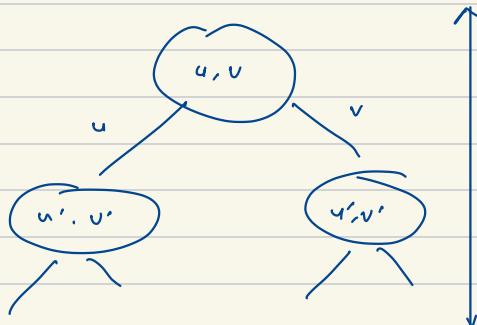
Guess - Try both

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

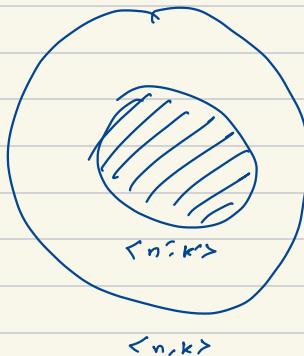
## Kernelization

$\rightarrow$  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

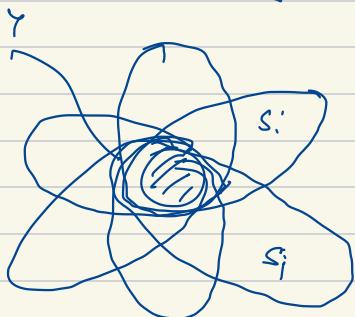
Sunflower is a collection of sets  $S_1, S_2, \dots, S_k$

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

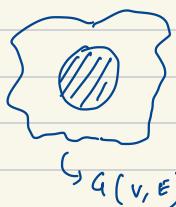
Fixed-parameter algo??

## 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

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

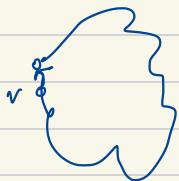
$\xrightarrow{\text{t}} \xrightarrow{\text{to success in one round}} \xrightarrow{\text{failure in one round}}$

# 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  $e^{kt+O(t)} / \gamma n$