



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

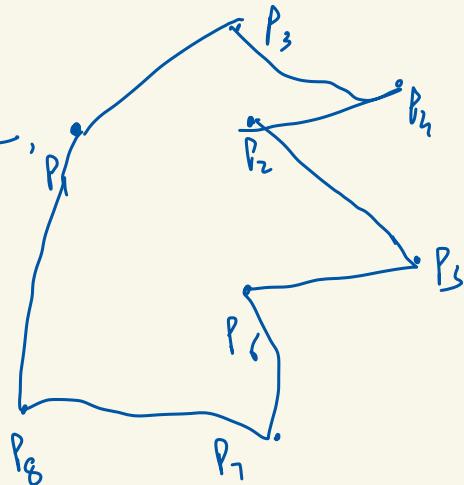
algorithm (of with some parameter $\epsilon > 0$) for any input, output a solⁿ within a factor $(1 + \epsilon)$ of the optimal solution.
 that runs in $n^{f(\frac{1}{\epsilon})}$ for some comparable f^n .
 (Running time is polynomial in $n, \text{some } \epsilon$)

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

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

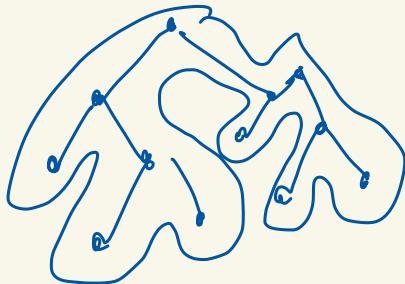
$S = \text{set of edges}$

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

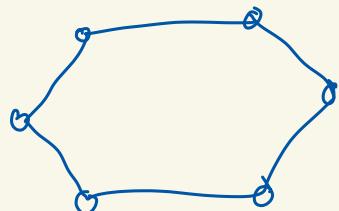
Every edge of the MST is traversed twice

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

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

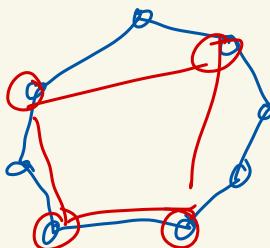
$$S \subseteq V$$

$$G_i(S)$$

$$\xrightarrow{\quad} OPT_S$$

$$OPT_S \leq OPT_G$$

Triangle
Inequality



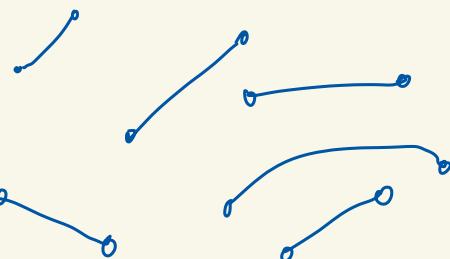
Property-2

Perfect Matching (can be computed in polynomial time)

Min cost AM

Perfect Matching with
smallest cost

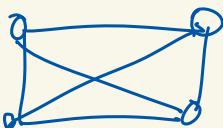
(Polynomial Time
 $\sim O(n^2)$)



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

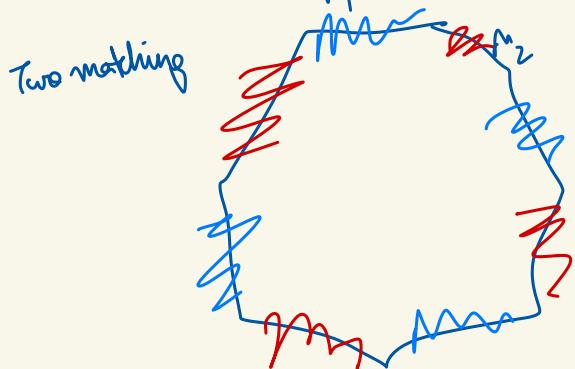
Compute Eulerian circuit

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

Analysis

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

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



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

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

$F = \{(i_1, j_1), \dots, (i_n, j_n)\}$ \rightarrow edges obtained

$OPT > \sum_{e=2}^n c_{i_e j_e}$

\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)

whereas

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

$$\text{cost (nearest-addition algo)} \leq 2 \sum_{e=2}^n c_{i_e j_e} \leq 2(OPT)$$

* 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 $(i_1, i_2), \dots, (i_{2k-1}, i_{2k})$

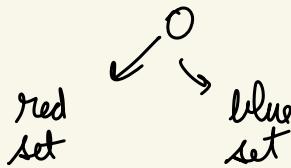
perfect-matching of minimum cost = $O(n^2)$

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

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

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

Consider the shortest tour on the node set \mathcal{D} . 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 \rightarrow 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 $\bar{s}(x) = \frac{|x \cap N|}{|N|}$

$Y \subseteq X$ $R_{1Y} = \{x \cap Y \mid x \in R\}$

projections of R on Y . The range space
is projected to Y is $S_{1Y} = (Y, R_{1Y})$

If R_{1Y} contains all subsets of Y ($\text{if } Y = \text{finite}, |R_{1Y}| = 2^{|Y|}$)

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

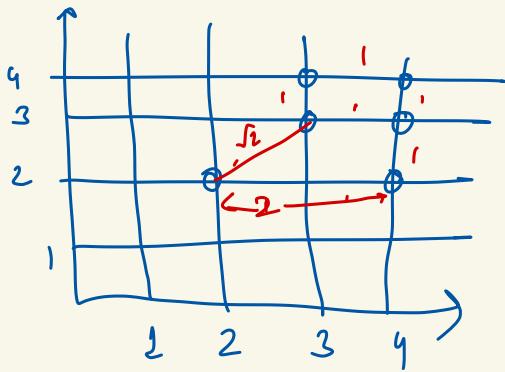
* Sauer's Lemma: If (X, R) is a range space of VC dimension S with $|X| = n$ then $|R| \leq G_S(n)$

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

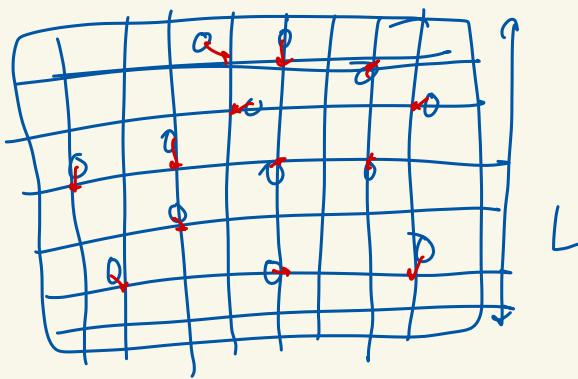
$$\{a_1, \dots, a_k\} \quad \{b_1, \dots, b_k\}$$

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

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

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

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

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

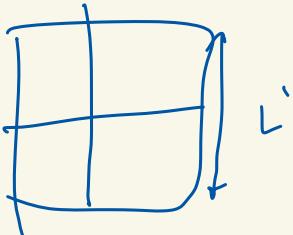
$$\text{OPT}_I + \epsilon \times \text{OPT}_I$$

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

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

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 → 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) $\leq (1+4\epsilon) \times \text{OPT}$.

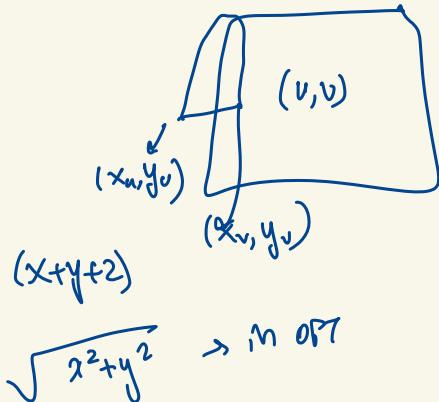
Extend non-ptour to a ptour

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 $(x+y+2)$



$$\sqrt{x^2 + y^2} \rightarrow \text{in opt}$$

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

Bound - expected length of the detours

Detours might when

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

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

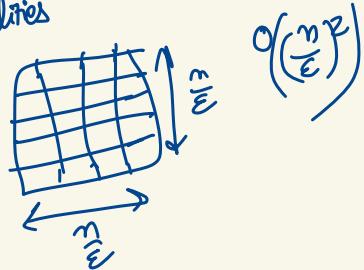
increase of detours $> 4\epsilon \text{ OPT}$)

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

De-randomize

- fixed ϵ

- grid shifting by trying all possibilities



Final Step (DP)

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

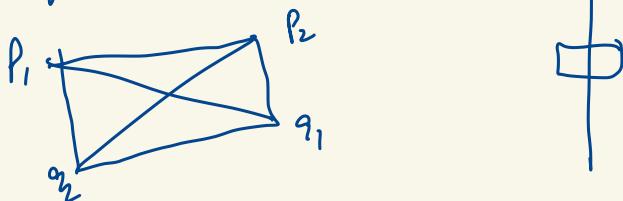
Introduce state

- 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

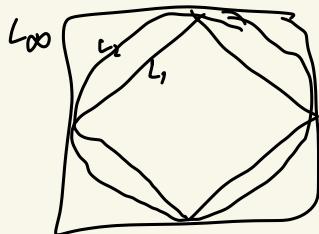
$p=1$: Manhattan distance (L_1 norm)

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 , $|P|=n$

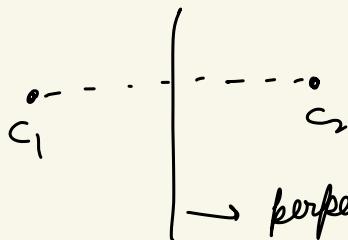
O/P : Find a clustering (set of centers) such that each pt is assigned to its nearest center

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



perpendicular bisector

now partitioned into two spaces

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

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 $\|P_K\|_1 \leq (1-\delta) \|P_{L_{\text{curr}}}\|_1$

- 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

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

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

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

K-means

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

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

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 $\tau \in R$

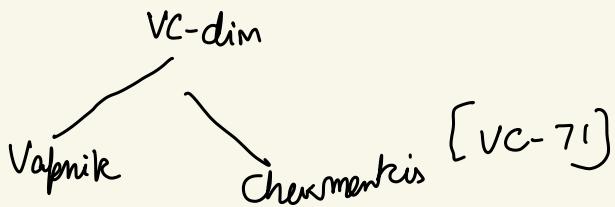
$$m(\tau) = \frac{|\tau \cap X|}{|X|}$$

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 Can 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

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

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

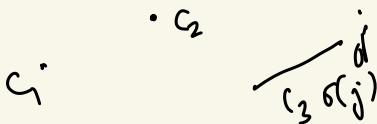
$$\overline{S} = (X, \overline{R}) \quad \text{where } \overline{R} = \{X \setminus \tau \mid \tau \in R\}$$

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)



$\text{opt} = X^*$ - optimal set of X $\leftarrow X - i + i'$
k-centers 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: Let $N \subseteq X$ 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

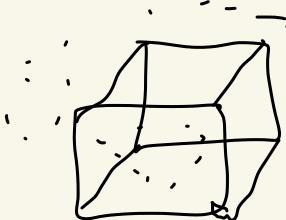
Allow 1% error,

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)

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

with prob $> 1 - \phi$.

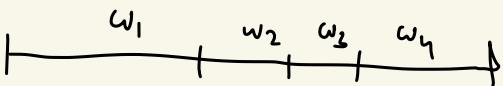
Suppose the shattering dim is d .

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

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. $\geq \epsilon w$



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^*

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

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

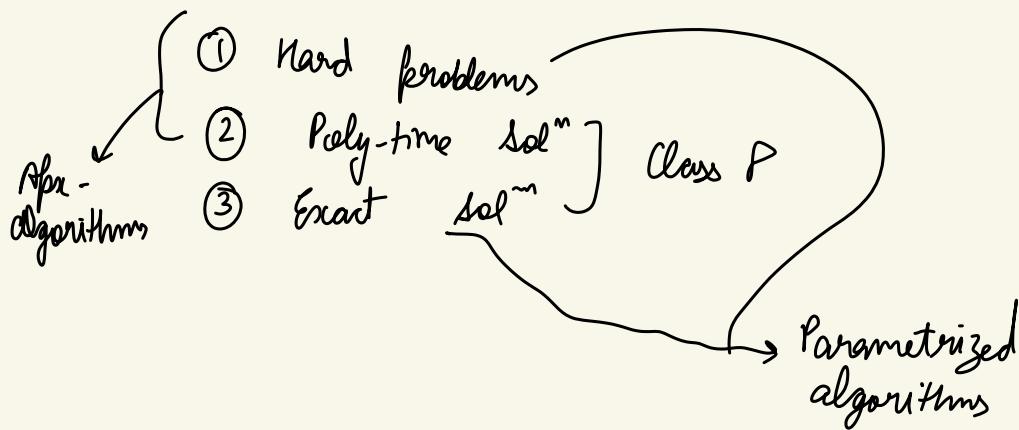
Q2 How to choose ε ?

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 solⁿ when the parameter small.

(Note: parameters are small in practice).

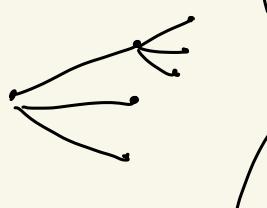
parameters - non-negative integer $k(x)$ (comes with prob i/p)
 - denote by k
 - Not necessarily efficiently computable.

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

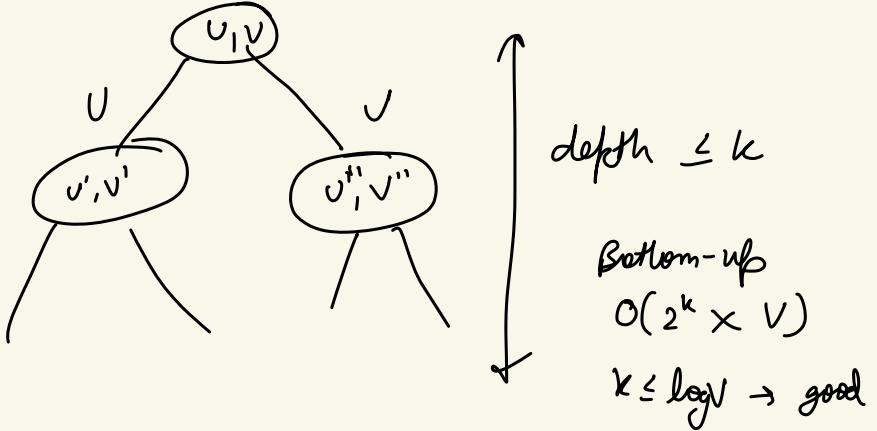
Guess — Try both

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 $f(k) \leq n$

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

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

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

Kernelization

simplifying
self-reduction

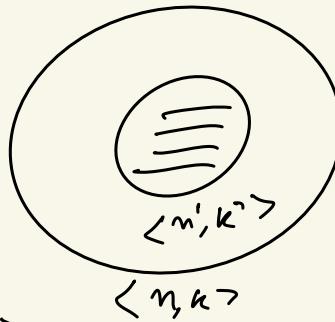
(poly-time reduction)

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)

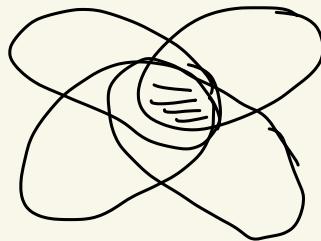
- Classical result from 1960

- 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 U , 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 $H \subseteq U$ 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 .

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

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