

Given 'n' objects, dissimilarity function
 $d(\cdot) : [n] \rightarrow \mathbb{R}_{\geq 0}$,
 group these into 'clusters' so that
 similar points are more likely in the
 same cluster.

Example Use Cases :-

- ① Categorizing data (websites/documents) based on content.
- ② Clustering songs into ragas.
- ③ cluster a city into neighborhoods to place utilities.
- ④ Ensuring diversity (first find clustering in a committee & then pick the committee by choosing from each cluster).

OUR MODELING in this COURSE

$\left. \begin{array}{l} G \\ V \end{array} \right\} - n \text{ points, in a generic 'metric space'}$
 way of formalizing
 a 'dis-similarity
 function'.

I
 V
 E
 N } - Target 'k'.
 } - objective function f

a dis-similarity function.

Goal: Partition $[n]$ to clusters $S_1 \uplus S_2 \uplus \dots \uplus S_k$

$$S_i \cap S_j = \emptyset$$

$$\bigcup_{i=1}^k S_i = [n]$$

to minimize $f(S_1, S_2, \dots, S_k)$

Find efficient (approximation) algo for this problem



WHAT IS A METRIC SPACE?

Given n points,
 a distance function $d(i, j)$ is
 a metric if it satisfies

- (i) $d(i, j) \geq 0 \quad \forall i, j \in [n]$
- (ii) $d(i, i) = 0 \quad \forall i \in [n]$

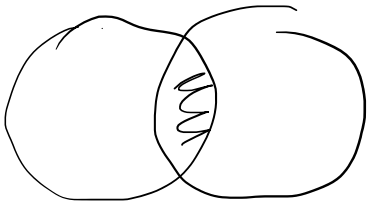
- (iii) $d(i, j) = d(j, i) \quad \forall i, j \in [n]$
 (iv) $d(i, j) \leq d(i, k) + d(k, j)$
 $\forall i, j, k \in [n].$

EXAMPLES OF METRIC SPACES :-

- ① - n points can be vectors in \mathbb{R}^d
 - $d(i, j) = \|v_i - v_j\|_2$
 or $\|v_i - v_j\|_p$ for any p .
 ✓ IS A METRIC
- ② n points can be vertices of $G = (V, E)$
 $d(i, j) =$ shortest path b/w i & j
 ✓ is a METRIC
- ③ If G is directed, it's not a metric.
 [Δ^k inequality is fine, but symmetry doesn't hold]
- ④ Have a collection of documents,
 each of which contains english words.
 let document i contain S_i (set of words)

then $J(i, j) = \text{JACCARD SIMILARITY}$

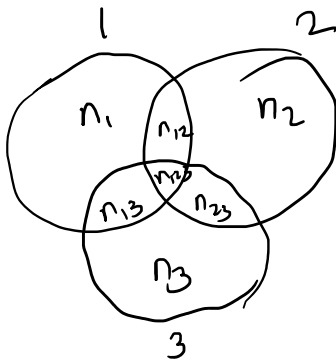
$$= \frac{|S_i \cap S_j|}{|S_i \cup S_j|}$$



then $d(i, j) = 1 - J(i, j)$

is a metric

(Δ^e inequality easy to see)



$$d(1, 2) = 1 - \frac{n_{12} + n_{123}}{T - n_3}$$

$$d(1, 3) = 1 - \frac{n_{13} + n_{123}}{T - n_2}$$

$$d(2, 3) = 1 - \frac{n_{23} + n_{123}}{T - n_1}$$

Need to show

$$d(1,2) \leq d(1,3) + d(2,3)$$

$$1 - \frac{n_2 + n_{23}}{T - n_3} \leq 1 - \frac{n_{13} + n_{123}}{T - n_2} + 1 - \frac{n_{23} + n_{123}}{T - n_1}$$

$$\frac{n_{13} + n_{123}}{T - n_2} + \frac{n_{23} + n_{123}}{T - n_1} \leq 1 + \frac{n_{12} + n_{123}}{T - n_3}$$

[It works,
check if true]

Common Objective Functions :-

① k-center Objective Function

(useful for placing police stations, etc)

First used in the 1950s, One of the earliest uses of approx Algs.

Basically, trying to min $\max_{k'=1}^k \max_{i,j \in S_{k'}} d(i,j)$

↓
(ie) min Max Diameter of each cluster.

Useful for Police Stations, School buses, Hospitals, etc.

FORMAL

Given a cluster S ,

$d(i,i)$

Given a cluster S ,

def. Center of cluster = $\arg \min_{i \in S} \max_{j \in S} d(i, j)$

* radius of cluster = $\min_{i \in S} \max_{j \in S} d(i, j)$

k-Center Problem

Given $[n]$ points, find a clustering of $[n]$ into S_1, S_2, \dots, S_k s.t we minimize maximum radius over all S_i .

↓

- Prefers all clusters are smallish over a clustering where one is very big & many are very small.

Possible Algorithms :-

① Can we use set cover?

Elements = $[n]$.

Sets correspond to

$$S_i = \{j : d(i, j) \leq r\}$$

for some 'r'. Then choose a min set cover.

- Keep increasing 'r' & stop when we can find a set cover of size 'k'.

② Idea: Obj fn. tries to avoid any point which is 'far' from all the k centers.

Greedy like Algo :-

① Pick C_1 as an arbitrary point from $[n]$.

② For $t = 2, \dots, k$

③ Pick $C_t = \operatorname{argmax}_{i \in [n]} \min_{1 \leq t' \leq t-1} d(i, C_{t'})$

④ Form the clusters with centers $\{C_1, C_2, \dots, C_k\}$ by assigning each point to nearest center

Idea: In step ③, we are checking if there is any point which hurts the obj. fn. for the current set of centers. We pick the 'worst' such point & make it a new cluster center.

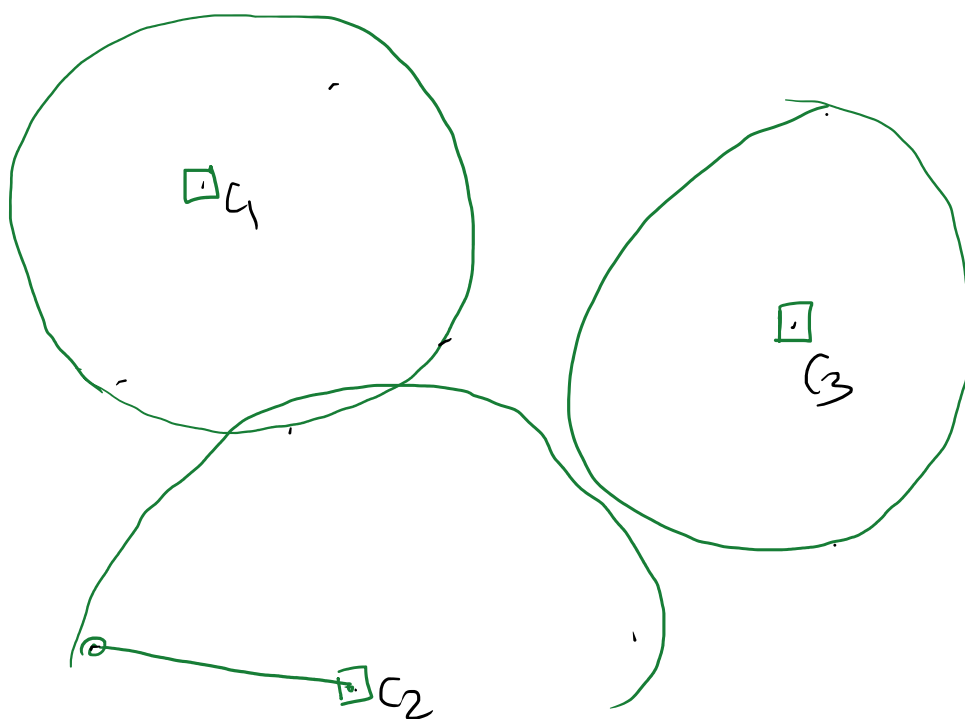
{ $\frac{AIM}{ALGO}$ is a 2-approx for k-center problem }

Given set P of n points in a metric space,
 choose a subset of k centers

$$C = \{c_1, c_2, \dots, c_k\} \subseteq P \text{ such that}$$

we can cluster the points in P with
 minimum radius clusters.

(ie) minimize $\max_{p \in P} \min_{c_i \in C} d(p, c_i)$.



Once we choose the centers, the
 clustering to minimize the max-radius
 is easy - each point is assigned to
 its nearest center.

Moreover, Max. radius of all clusters

$$= \max_{P \in P} \min_{C_i \in C} d(P, C_i)$$

let's focus on the distance of each point to its nearest center

once we fix the k centers
dist-to-closest-center

P_1	5
P_2	10
P_3	12
\vdots	2
\vdots	0
P_n	7

Restating the problem:

Choose k centers $\{C_1, C_2, \dots, C_k\} \subseteq P$

to minimize

$$\max_{P \in P} \min_{C_i \in C} d(P, C_i)$$

— Motivates the greedy algorithm

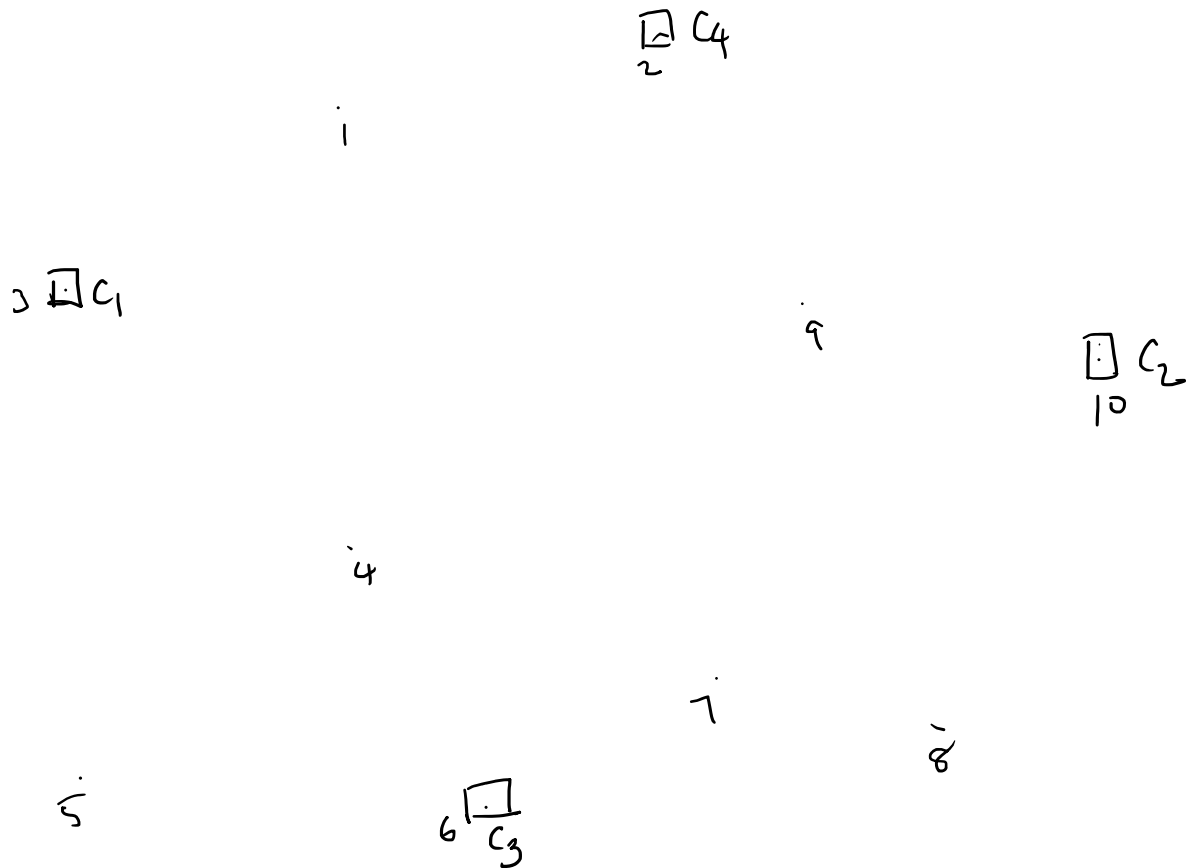
① choose c_1 arbitrarily from P

② for $t = 2, \dots, k$

③ let $c_t = \underset{p \in P}{\operatorname{argmax}} \min_{1 \leq t' \leq t-1} d(p, c_{t'})$.

④ Form the clustering using c_1, c_2, \dots, c_k as centers. \square

Toy Example



Suppose $k = 4$

THEOREM

*

If optimal clustering has objective value R ,
 our clustering has objective value $\leq 2R^*$.
 (\Rightarrow 2-approximation).

Proof:-

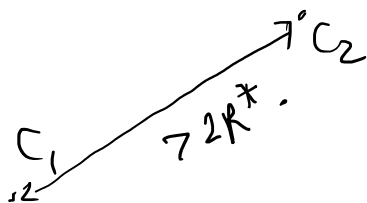
Let c_1, c_2, \dots, c_k be algorithm's centers.

Let $c_1^*, c_2^*, \dots, c_k^*$ be optimal centers.

Suppose for contradiction, our clustering
 has objective value $> 2R^*$.

$\Rightarrow \exists$ a point $\hat{p} \in P$ st it is
 far from all c_1, c_2, \dots, c_k

$\Rightarrow d(\hat{p}, c_i) > 2R^* \quad \forall i = 1, 2, \dots, k.$



c_4

c_5

c_3

\hat{p}

* \hat{p}

Because greedy also chose

$$C_t = \operatorname{argmax}_{P \in \mathcal{P}} \min_{1 \leq t' \leq t-1} d(P, C_{t'})$$

and $\min_{1 \leq t' \leq t} d(\hat{p}, C_{t'}) > 2R^*$

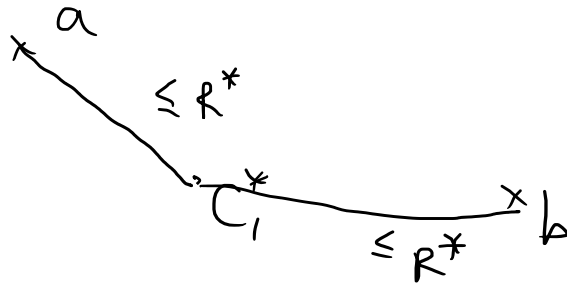
$\Rightarrow \hat{p}$ is a candidate for all of C_1, \dots, C_k but didn't get picked

$\Rightarrow \forall t, \min_{1 \leq t' \leq t-1} d(C_t, C_{t'}) > 2R^*$ as well.

\Rightarrow There exist $k+1$ points in \mathcal{P} namely $\{C_1, C_2, \dots, C_k, \hat{p}\}$ st each pair of points is at $> 2R^*$ distance.

But now how are these points clustered well in OPT?

\exists some cluster, say C_1^* with
 ≥ 2 points a and b
 from $\{C_1, C_2, \dots, C_k, \hat{V}\}$



Because OPT has radius R^*

$$d(a, C_1^*) \leq R^* \quad \text{and}$$

$$d(b, C_1^*) \leq R^*$$

Now, Δ^e inequality gives

$$d(a, b) \leq d(a, C_1^*) + d(C_1^*, b) \leq 2R^*$$

[CONTRADICTION]

\square

Same type of input as k-center, but a slightly different objective function.

Given n points in a metric space

Find k centers c_1, c_2, \dots, c_k

to minimize

$$\sum_{P \in P} \min_{i=1}^k d(P, c_i)$$

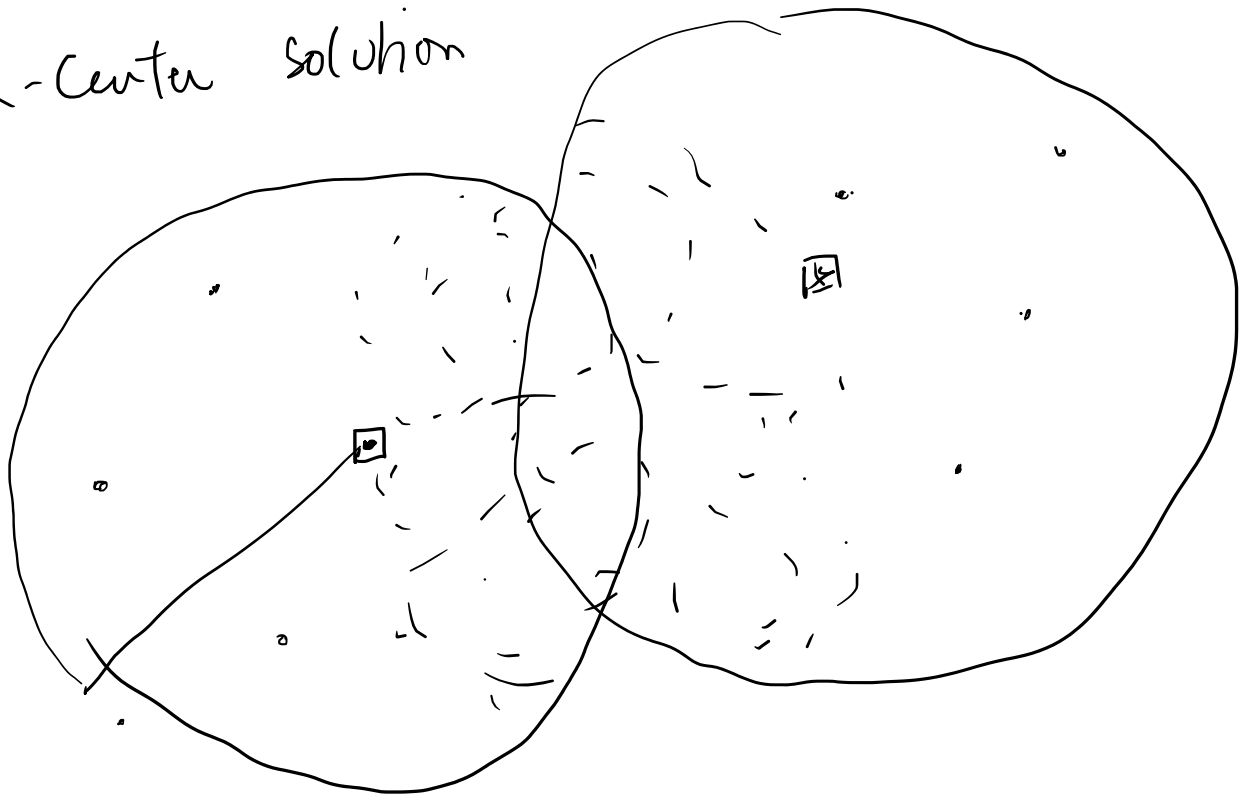
Motivation of Objective Function

① Maybe we are trying to lay cables from k stations to all points of the city, where to place the power stations?

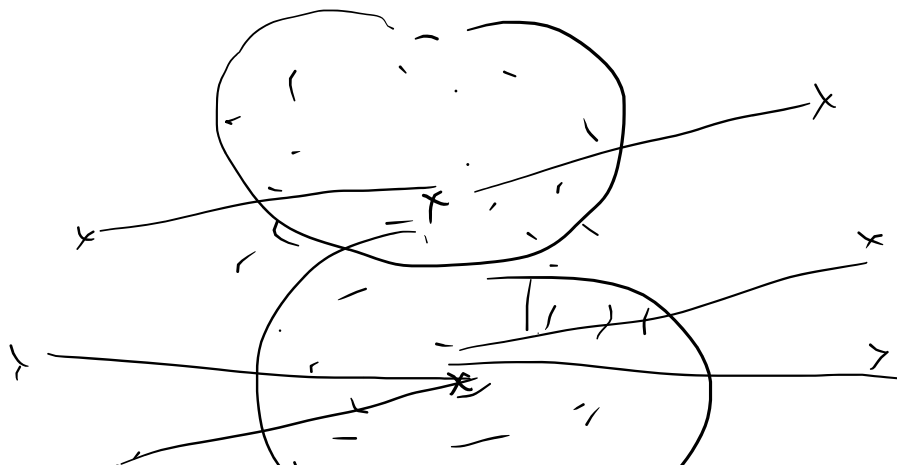
Total cable length matters more than max.



k-Center solution



k-Median solution





Try to cover the dense region better while being OK with a few points paying a large cost.

Yet another problem

k-Means Problem.

Same as above, objective is

$$\text{minimize} \sum_{P \in P} \min_{i=1}^k d(P, C_i)^2$$

Has lots of applications in ML, especially when points are vectors in \mathbb{R}^d and

$$d(i, j) = \|v_i - v_j\|_2$$

Has very nice physics connections to concepts such as

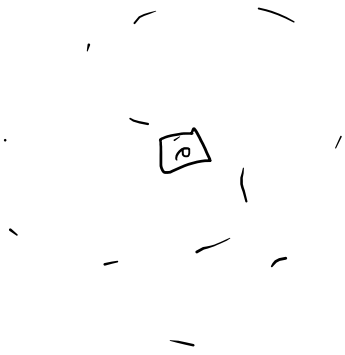
Center of gravity, etc.



Point which minimizes
sum of squared
distances is the
'mean' / centroid
of dataset.



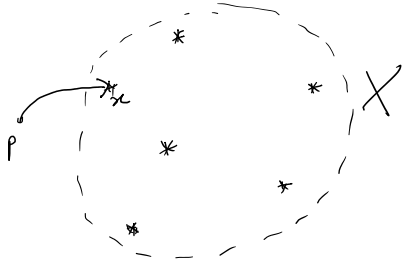
k-Means is a natural generalization
to k-clusters.



Given n points in a metric space P
 find k centers $G = \{c_1, c_2, \dots, c_k\}$

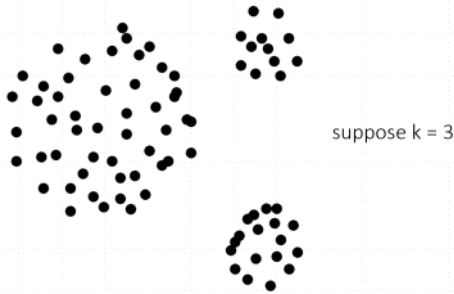
to minimize $\sum_{p \in P} d(p, G)$

where $d(p, X) = \min_{x \in X} d(p, x)$

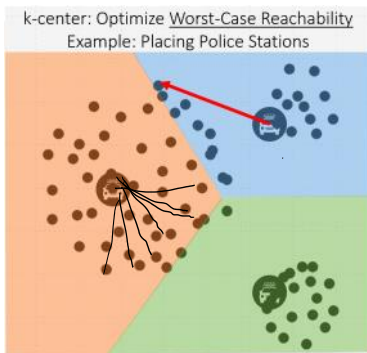


Recall k -Center objective was to
 minimize $\max_{p \in P} d(p, G)$

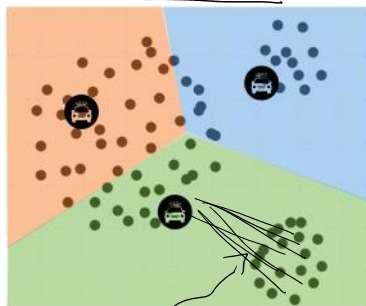
Illustration of k -Median vs k -Center :-



k -Center Solution



k -Median Solution



Maybe its ok for a few points to pay a large cost if many points pay less.

k -Center admits 2-approximation using greedy algorithm. What about k -Median?

Greedy algorithm. What about k-Median?

Simple greedy-type Alg's don't end-up being very good, we can resort to Linear Programming.

Y_i = variable for whether i is chosen as a center or not.

X_{ij} = variable for whether point j is assigned / clustered to center i .

LINEAR PROGRAM (K-MEDIAN)

$$\text{Min } \sum_j \left(\sum_i d(i,j) \cdot x_{ij} \right)$$

$$\sum_{i=1}^k x_{ij} \geq 1 \quad \forall j \in P$$

$$\sum_{i=1}^k Y_i \leq k$$

$$x_{ij} \leq Y_i \quad \forall i, j$$

$$x_{ij} \geq 0$$

$$Y_i \geq 0$$

Lemma 1

Let (x^*, y^*) be an optimal LP solution

$$\text{Then } LP^* = \sum_j \sum_i d(i,j) x_{ij}^* \leq OPT$$

where OPT is the k-Median cost of optimal solution.

Proof

Unknown optimal k-Median solution is feasible for the LP, which can only do better.

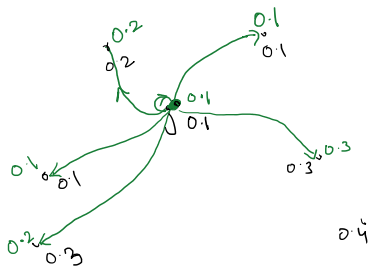
QUESTION

How do we "round" this fractional solution into a good clustering?

Also note :-

In the optimal LP solution, once we know the y^* values, the x_{ij}^* values can be easily derived

Let's look at some point j .



White values are y_i^* values.
Green values are x_{ij}^* values.

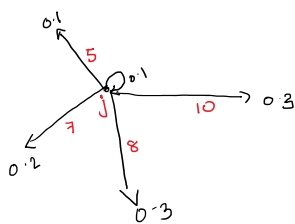
Idea

infer some basic properties of what the LP optimal is trying to do.

In particular, let

$D_j =$ LP-distance that point j incurs in the optimal solⁿ

$$= \sum_i d(i,j) \cdot x_{ij}^*$$



→ Distances are in Red

→ x_{ij}^* values are in white

$$\begin{aligned} D_j &= \text{LP-distance of point } j \\ &= 0.1 \times 5 + 0.2 \times 7 + 0.3 \times 10 + 0.4 \times 8 \\ &= 0.5 + 1.4 + 3 + 3.2 \\ &= 7.3 \end{aligned}$$

LP tries to cover point j by connecting it to a center at "distance" 7.3, so we use that as a guide.

LEMMA 2

For any point j , let $B_j = \{i : d(i, j) \leq 2D_j\}$
 be the points at distance $\leq 2D_j$
 from j .

Then $\sum_{i \in B_j} y_i^* \geq \sum_{i \in B_j} x_{ij}^* \geq \frac{1}{2}$

Proof :-

Suppose not, and $\exists j$ st

$$\sum_{i \in B_j} x_{ij}^* < \frac{1}{2}$$

Then

$$\begin{aligned} D_j &= \sum_i d(i, j) x_{ij}^* = \sum_{i \in B_j} x_{ij}^* d(i, j) + \\ &\quad \underbrace{\sum_{i \notin B_j} x_{ij}^* d(i, j)}_{> 0} \\ &> 0 + \left(\sum_{i \in B_j} x_{ij}^* \right) \cdot 2D_j \\ &> 0 + \frac{1}{2} \cdot 2D_j \\ &= D_j \quad \Rightarrow \Leftarrow \end{aligned}$$

for each j , if we place some center
 at a point in B_j , then its
 connection distance $\leq 2D_j$

$$\Rightarrow \text{Overall cost} \leq \sum_j 2D_j \leq 2LP^* \leq 2OPT$$

12/03/2021

Algo 1

① Sort j such that

$$D_{j_1} \leq D_{j_2} \leq \dots \leq D_{j_n}$$

② Pick the set "near-independent"
 points J^* as follows

- For $l = 1 \dots n$ (in sorted order)

if there exists no $j' \in J^*$

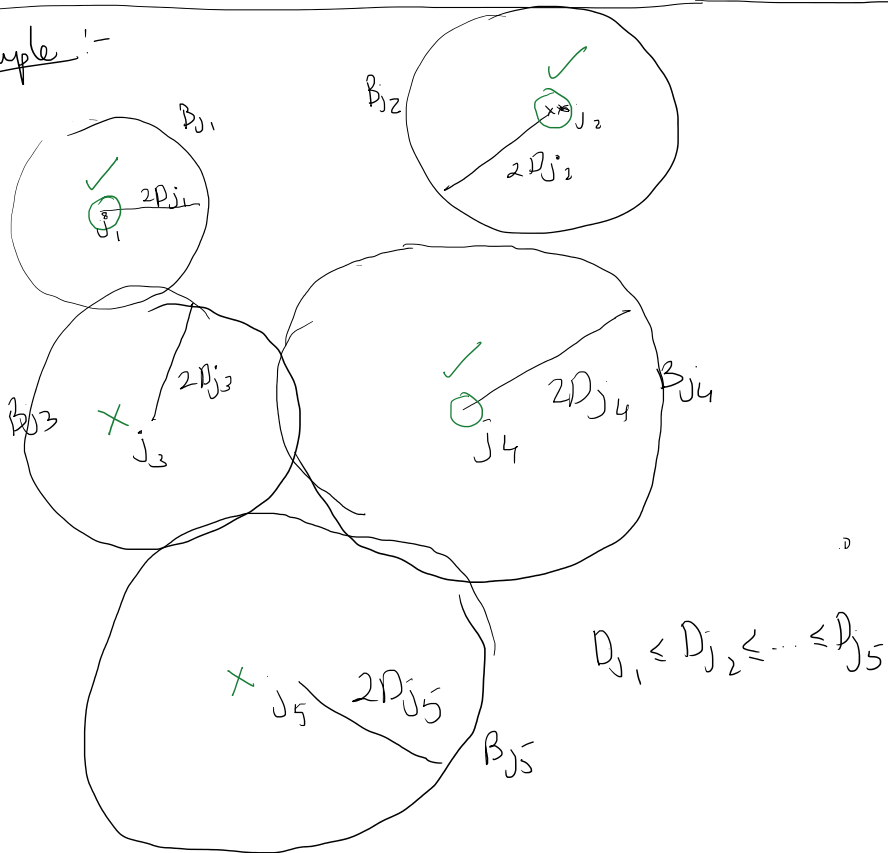
$$\text{st } d(j_l, j') \leq 2D_{j_l} + 2D_{j'}$$

then add j_l to J^* .

...

③ Open a center at each $j \in J^*$

Example :-



Finally, J^* contains $j_1, j_2,$ and j_4 .

Lemma ①

J^* is s.t. $\forall j', j'' \in J^*$
 $B_{j'} \cap B_{j''} = \emptyset$

Proof

Sup $B_{j'} \cap B_{j''} \neq \emptyset$ and say i belongs to both.

$$\begin{aligned} \text{then } d(j', j'') &\leq d(j', i) + d(i, j'') \\ &\leq 2D_{j'} + 2D_{j''} \end{aligned}$$

(contradiction to whicher got added later)

What did we do?

① Somehow identify points which are

not too overlapping

② Remaining points are "close" to the J^* points \Rightarrow if we can handle the J^* , we can hope to handle the rest also.

SIMPLE ALGO ① { Won't exactly give k centers, but might open up to $2k$ centers }

For each $j \in J^*$, open a center at j .

Claim ①

$\forall j$, distance of j to nearest open center $\leq 4D_j$

Claim ②

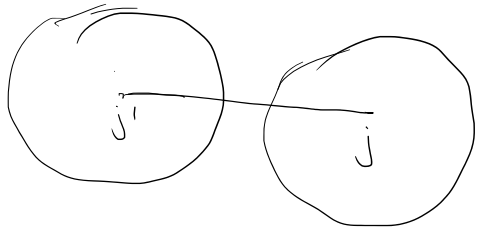
Total # open centers $\leq 2k$.

Proof

Claim ① :-

If j got added to J^* , then dist of j to nearest center = 0 $\leq D_j$ ✓

If j didn't get added, there must exist a j' st



$$d(j, j') \leq 2D_j + 2D_{j'} \\ \leq 4D_j \quad \checkmark$$

So dist to nearest open center $\leq 4D_j$ 😊

Pf of claim ② →

In each B_j , $\sum_{i \in B_j} y_i^* \geq \frac{1}{2}$

Moreover $B_{j_1} \cap B_{j_2} = \emptyset$ for $j_1, j_2 \in J^*$

So, simply sum over all $j \in J^*$

$$\sum_{j \in J^*} \sum_{i \in B_j} y_i^* \leq \sum_{i=1}^n y_i^* \leq k \quad \uparrow \text{LF constraint}$$

$$\sum_{j \in J^*} \sum_{i \in B_j} y_i^* \geq \sum_{j \in J^*} \frac{1}{2} = \frac{|J^*|}{2}$$

$$\Rightarrow |J^*| \leq 2k$$

Bi-Criteria Approximation Algorithm

Given an instance of k -Median, with optimal cost $= \text{OPT}$, we can efficiently find a solution which opens $2k$ centers and has cost $\leq 4 \text{OPT}$

$\leq 4 \text{OPT}$

How do we improve to a pure
k-Median solution?

I
D
E
A

Focus on J^* , pair them up in

J^* , so that

in each pair $\sum y_i^* \geq 1$

and then handle each pair

Recap

- 1) Solve LP
- 2) Get (x^*, y^*) as solution
- 3) Define $D_j = \sum_i d(i, j) x_i^*$
- 4) Define $B_j = \{i : d(i, j) \leq 2D_j\}$
- 5) $y^*(B_j) = \sum_{i \in B_j} y_i^* \geq \frac{1}{2} \quad \forall j$
- 6) J^* is the near-independent far-away points in J .
- 7) $\forall j \in J^*, \exists j' \in J^*$ st $d(j, j') \leq 4D_j$
and $D_{j'} \leq D_j$
- 8) $|J^*| \leq 2k$
- 9) $\forall j, j' \in J^*, B_j \cap B_{j'} = \emptyset$

↓
 From 1-8, we got a 4-approximation

- ① which opens $(2k)$ cluster centers
- ② Today, we'll try to make it a genuine k -clustering which opens $(\leq k)$ centers while being slightly worse in cost.

Idea

- [Pair up points in J^*
- [Open one center in each pair]

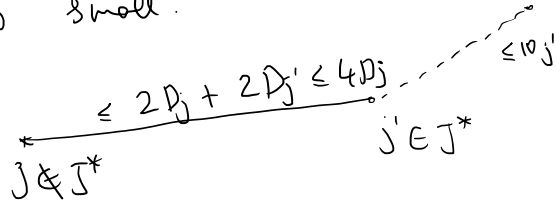
$|J^*| \leq 2k \Rightarrow \# \text{ centers open } \leq k. \text{ 😊}$

- How do we analyze the cost?
- How to pair points in J^* ?
- How to decide which center to open in a pair?

a pair ?

We'll try to ensure the connection cost of each point $j \in J^*$ is at most, say, $10 D_j$.

\Rightarrow Connection cost of points not in J^* is also small.



$$\begin{aligned} \Rightarrow \text{conn. cost of } j &\leq \\ &2D_j + 2D_{j'} + 10D_{j'} \\ &\leq 14D_j \end{aligned}$$

\Rightarrow It suffices to show that points in J^* have low connection cost.

Q How do we choose k centers to ensure that points in J^* have low conn cost ($\propto D_j$)

PAIRING ALGORITHM

- Among J^* , choose the closest pair of points say (j_1, j_2) and match them.
- Remove j_1, j_2 from J^* and repeat

$$\uparrow \text{ \# pairs + (singletons, if left) } \leq k$$

For now, let's assume no singleton left.
 \downarrow can handle it easily later.

↓
can handle it easily later.

We'll come up with a randomized selection procedure which ensures

- ① Each point is chosen as a center with prob y_i^*
- ② Total # of centers opened $\leq k$
- ③ For each matched pair (j_1, j_2) , we definitely open ≥ 1 center from among the points $B_{j_1} \cup B_{j_2}$.

SATISFYING just ① is easy.

Each i will independently choose itself as a center w.p. y_i^* .

↓
As a solⁿ idea, not great because there could be a region in space where we don't open any center



Independent rounding

$$\Rightarrow (1 - 0.1)^{10} \approx \frac{1}{6} \text{ prob.}$$

not opening any center in this cloud

Very few!!

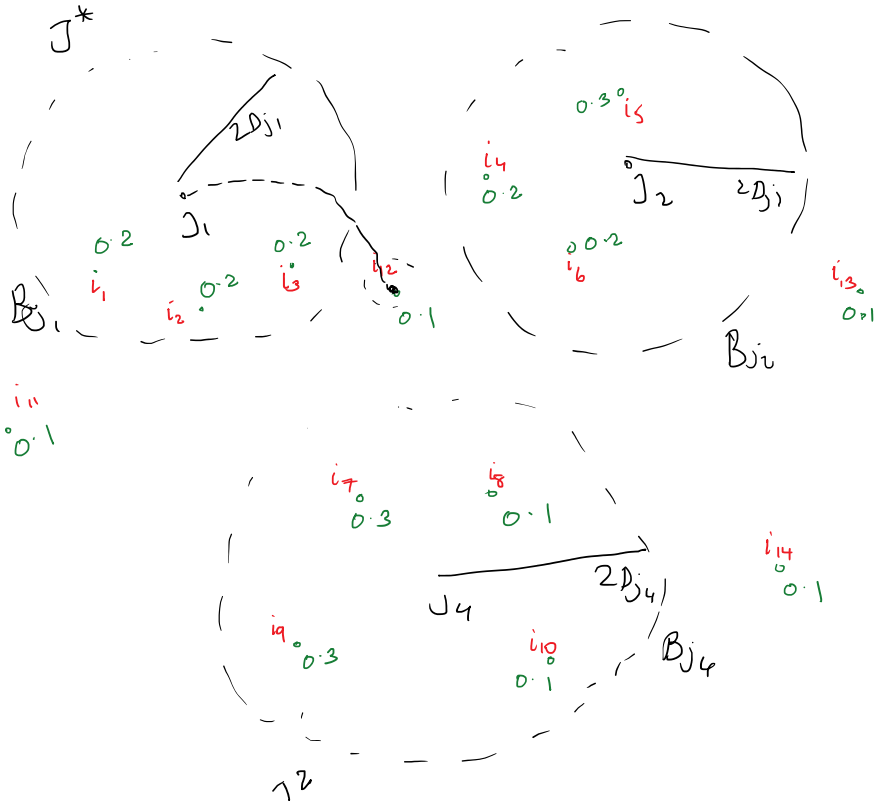
It fails to satisfy ② & ③ also.

In fact ③ precisely tries to address the issue of ^{completely} missing local regions in space, which is the drawback of independent rounding.

↓
 α -point rounding / dependant rounding
- Nice properties of prob αy_i^*

along with satisfying extra constraints.

Selection Procedure Illustration



Example in picture above
 J_1 connects to $i_1, i_2, i_3, i_{12}, i_{11}, i_4$

If i_i , we want to associate y_i length segment on the real line, and put them consecutively

Attempt ①: Order arbitrarily i_1, i_2, \dots, i_n



and place segments consecutively

Pick random $\alpha \in [0, 1]$ uniformly and
 Mark $\alpha, 1+\alpha, 2+\alpha, \dots, k-1+\alpha$

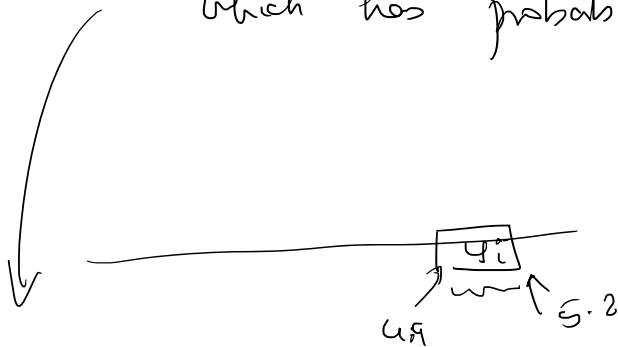
Because $\sum y_i = k$, length of line
segment $l = k$.

If a dart intersects a segment
 y_i , open a center at i .

Nice properties:-

$\forall i$, $\Pr [i \text{ is opened as center}]$

First dart crosses i , only if $0 \leq x \leq y_i$,
which has probability (y_i)



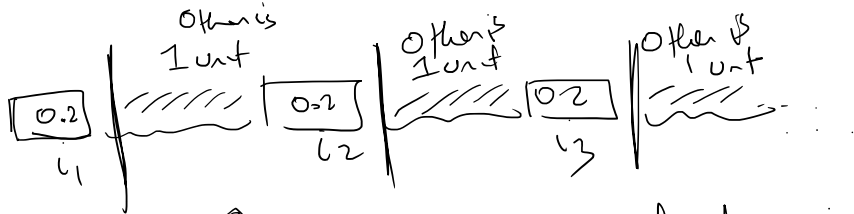
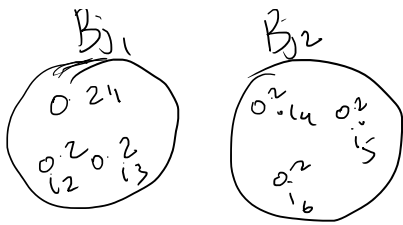
This is true for any i ,

$$\Pr [i \text{ is chosen as center}] = y_i$$

② ✓ # centers = k
because $(k+1)^{\text{th}}$ dart will lie
outside the system.

③ X doesn't satisfy the property
that in a "local region",
there is one open center

(e.g) $\forall (i, j) \in M$, we may not
open any center in $B_{i_1} \cup B_{i_2}$
 B_{i_1} B_{i_2}



↑ if $\alpha > 0.2$, all dark miss B_{j_1} and B_{j_2}

So, need to preserve locality in some manner.

Goals:

- ① Ensure B_j is contiguous for each j
- ② Ensure B_{j_1} and B_{j_2} are contiguous for each $(j_1, j_2) \in M$
- ③ Ensure other close points to j are contiguous for each j

for each i not in any B_j for $j \in J^*$,
 ensure it is contiguous with its nearest $j \in J^*$

$$O_j = \left\{ i : d(i, j) < d(i, j') \ \forall j' \in J^* \atop j' \neq j \right\}$$

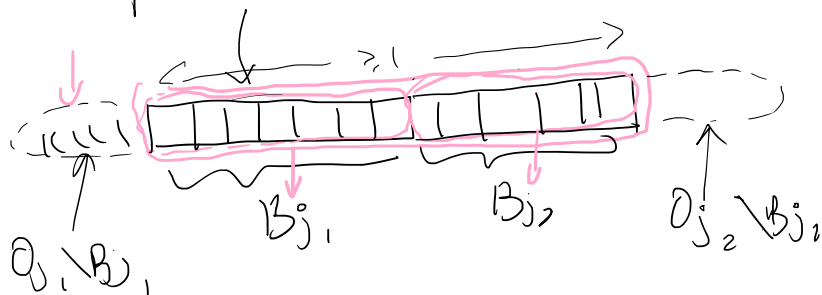
No ties in distances } lets assume that we break ties arbitrarily.

no, we can } we break it arbitrarily.
 ↓ (1) $d_1 \neq d_2$ for all pairs

(3) $\forall i, j$, want O_j appears consecutively

Will this work?

Take pair $(i_1, i_2) \in M$



repeat for all pairs

↓

Now since $y(B_{j_1}) + y(B_{j_2}) \geq 1$

we will definitely open
 one center in each "pair".

with same left throwing algorithm.

17/03/2021

Algorithm

- Recall defn of $B_j, O_j, T^*, D_j, x^*, y^*$
- Recall Pairing M .

- Place the points on a line and do
 the "α-point rounding"

↓
 Ensure that for each pair,

B_{j_1} is contiguous

B_{j_2} is contiguous

$O_{j_1} \setminus B_{j_1}$ contiguous with B_{j_1}

$O_{j_2} \setminus B_{j_2}$ contiguous with B_{j_2}

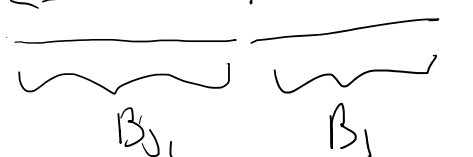
- Open all centers which are crossed by the darts

Lemma ①

Each i is opened with prob y_i^*

Lemma ②

For any pair $(i, j) \in M$, at least one point is opened from $B_{j_1} \cup B_{j_2}$ with probability 1.

$\leftarrow \text{length} \geq 1 \rightarrow$ because $y^*(B_{j_1}) + y^*(B_{j_2}) \geq 1$

and $B_{j_1} \cap B_{j_2} = \emptyset$



gap b/w any 2 darts = 1

ENSURES "LOCALITY PRESERVING ROUNDING"

Lemma ③

For each point $\bar{j} \in J^*$, let

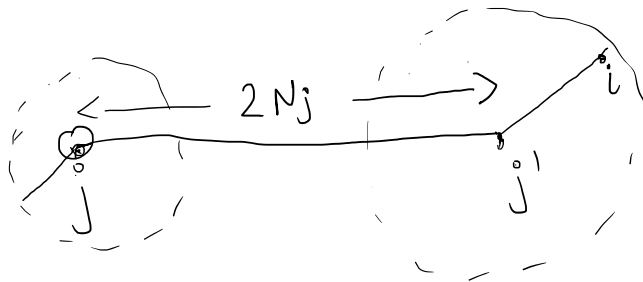
$N_{\bar{j}} = \frac{1}{2} \min_{i \in J^*} d(\bar{j}, i)$ be half

$$N_j^o = \frac{1}{2} \min_{\substack{j' \in J^* \\ j' \neq j}} d(j, j') \quad \text{be half the dist. to nearest other pt in } J^*$$

Then, there is always an open center within $6 N_j^o$

Proof

Let's focus on j and j' - its nearest other pt from J^*



Either $(j, j') \in M$ or not

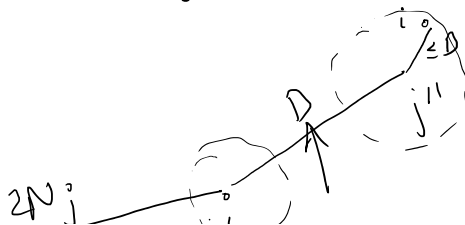
Case ①: If $(j, j') \in M$.

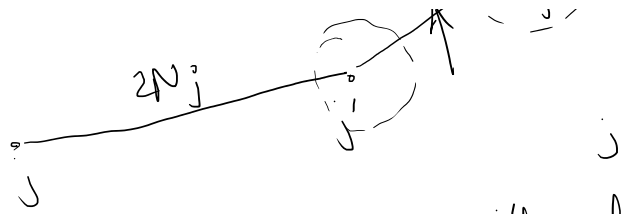
There is open center among $B_j \cup B_{j'}$

$$\Rightarrow d(j, i) \leq 4N_j^o$$

Case ② $(j, j') \notin M$.

$\Rightarrow \exists j''$ st $d(j', j'') \leq 2N_j^o$ and $(j', j'') \in M$.





But again, there will always be an open center in $B_j \cup B_j''$

$$d(j, i) \leq 2N_j + D + D$$

$$\leq 6N_j \quad (\Delta^k \text{ inequality})$$

Rest of Analysis :- \leftarrow LP connection cost.

If for some j , $D_j \geq 0.1N_j$ then

such j 's are very happy

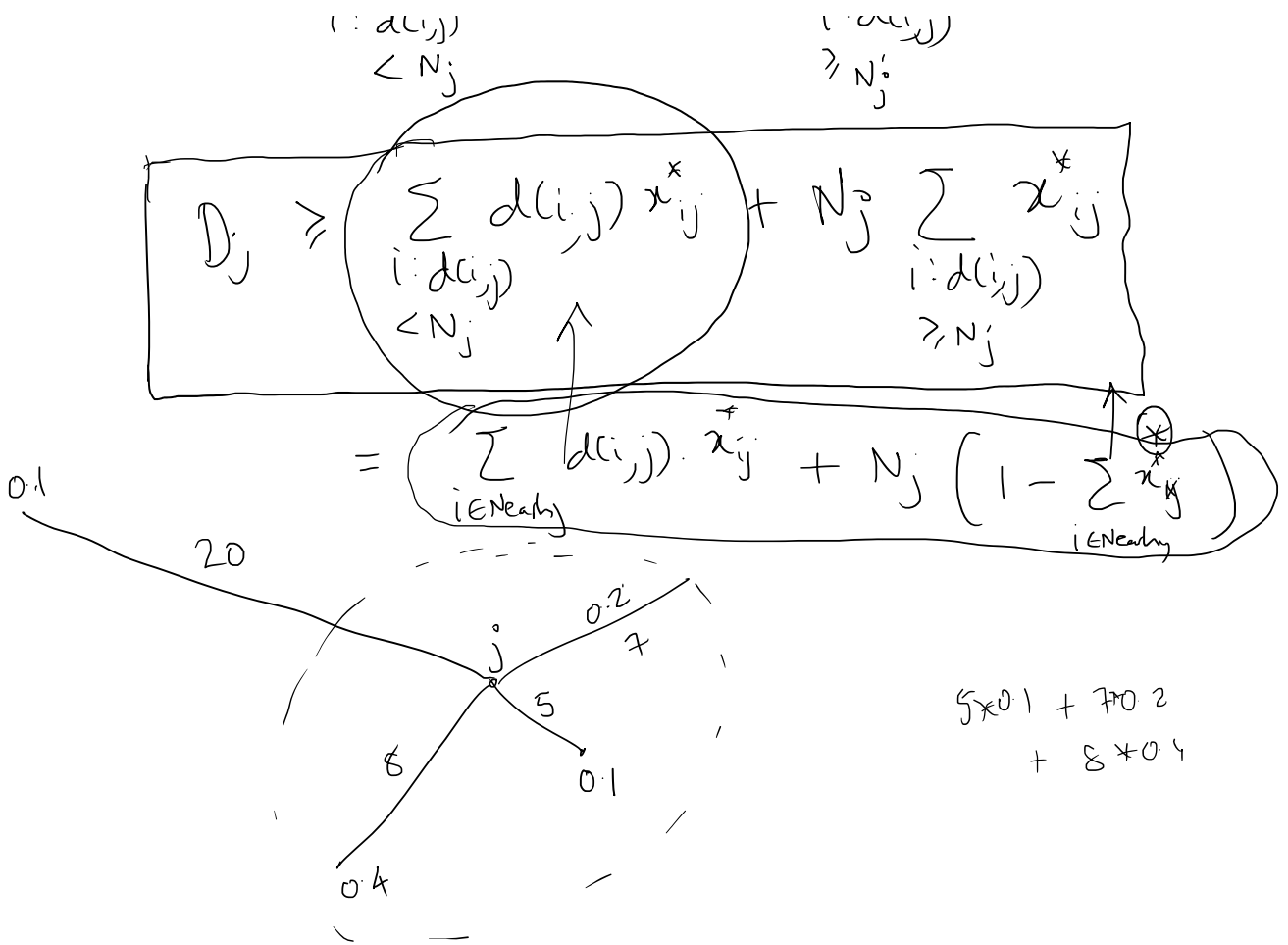
(Their real connection cost ≤ 60 their LP connection cost)

Real problem happens if D_j is much smaller than N_j

Let's analyze D_j^* :-

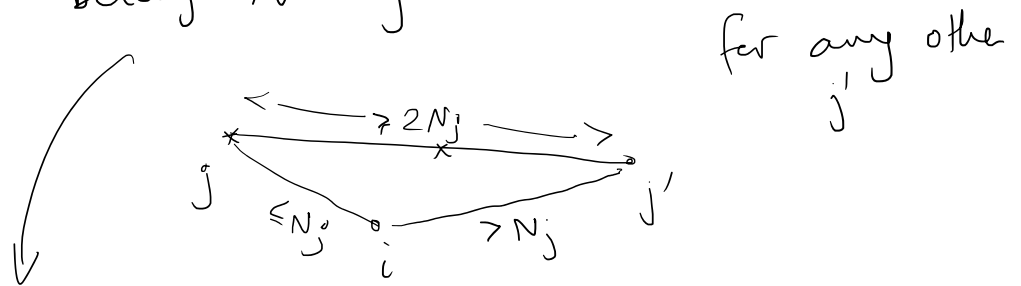
$$D_j^* = \sum_i d(i, j) x_{ij}^*$$

$$= \sum_{\substack{i: d(i, j) \\ < N_j}} d(i, j) x_{ij}^* + \sum_{\substack{i: d(i, j) \\ \geq N_j}} d(i, j) x_{ij}^*$$



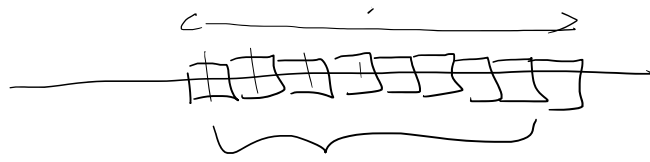
What's the expected connection cost of this point j , in our rounding scheme?

Notice: all the points i st $d(i,j) < N_j$ belong to O_j



Now, we ensured that such points can be placed contiguously in the line.

All points $i \in \{d(i, j) < N_j\}$ are placed contiguously.



$i : d(i, j) < N_j \leftarrow$ Nearby points

Expected cost of j 's connection

$$= \Pr(\text{one of the nearby pts chosen}) \cdot E[\text{distance} \mid \text{nearby points chosen}]$$

$$+ \Pr(\text{Nearby pt not chosen}) \cdot E[\text{distance} \mid \text{Nearby pt not chosen}]$$

$$= \left(\sum_{i: d(i, j) < N_j} y_i^* \right) \left[\frac{\sum_{i: d(i, j) < N_j} d_{ij} y_i^*}{\sum_{i: d(i, j) < N_j} y_i^*} \right] + \left(1 - \sum_{i: d(i, j) < N_j} y_i^* \right) \cdot 6N_j$$

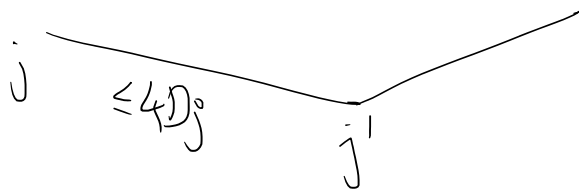
$$\leq 6D_j \quad \text{😊}$$

Corollary

$$\forall j \in J^*, E[\text{connection cost of } j] \leq 6D_j$$

Corollary

$$\forall j \in J^*, E[\text{conn. cost of } j] \leq 10D_j$$



$$\exists j' \text{ st } d(j, j') \leq 4D_j \ \&$$

$$D_{j'} \leq D_j$$

$$\Rightarrow E[\text{conn. cost of } j] \leq 4D_j + E[\text{conn. cost of } j']$$

$$\leq 4D_j + 6D_j$$

$$\leq 10D_j$$

Takeaway

① $O(n)$ for k-Median

② Respecting hard constraints $\sum y_i \leq k$ during rounding is hard, we need "dependent rounding" ideas

③ α -point rounding is a good way to ensure this.

Next 2-3 lectures, yet another algo. for k-Median which gets rid of ② in a clever way.

Given metric (X, d) $|X| = n$ points

$d(\cdot)$ is distance function (metric)

- $d(i, j) + d(j, k) \geq d(i, k) \quad \forall i, j, k \in [n]$
- $d(i, j) = d(j, i)$
- $d(i, i) = 0$

Choose k points as centers and assign each point in X to nearest center to minimize total "assignment distance"

(ie) $\sum_{j=1}^n d(j, S)$ where $d(j, S) = \min_{i \in S} d(i, j)$

and $|S| = k$

Recall LP

Total Assignment distance \leftarrow Min $\sum_j \sum_i d(i, j) x_{ij}$

every point is assigned \leftarrow $\sum_i x_{ij} \geq 1 \quad \forall j$

center must be open \leftarrow $x_{ij} \leq y_i \quad \forall i, j$

$\sum y_i \leq k$

$$\left. \begin{array}{l} \text{Total } k \text{ centers} \leftarrow \sum_j d_{ij} = 1 \\ x_{ij} \geq 0 \\ y_i \geq 0 \end{array} \right\}$$

Method of Lagrangian Relaxation

Idea:

We had a lot of difficulty in rounding the LP solution to preserve $\sum y_i \leq k$ (allowing violations was much easier)

Instead, let us push this constraint to the objective function

New LP: LP2

$$\text{Min}_{x,y} \sum_j \sum_i d_{ij} x_{ij} + \lambda (\underbrace{\sum y_i - k}_{\text{constraint}})$$

$$\text{LP2} \left\{ \begin{array}{l} \sum_i x_{ij} \geq 1 \quad \forall j \\ x_{ij} \leq y_i \quad \forall i, j \\ x_{ij} \geq 0 \\ y_i \geq 0 \end{array} \right\} \forall i, j$$

$$\forall \lambda \geq 0,$$

$$\text{OPT}(\text{LP2}) \leq \text{OPT}(\text{LP}).$$

In LP2, we can effectively ignore the $-\lambda k$ constant in the obj (It is the same for all x, y solutions).

LP3 is LP2 without $-\lambda k$ term

$$\begin{aligned} \text{LP3} \quad \text{Min} \quad & \sum_j \sum_i d(i,j) x_{ij} + \lambda \sum_i y_i \\ & \sum_i x_{ij} \geq 1 \quad \forall j \\ & x_{ij} \leq y_i \quad \forall i, j \\ & x_{ij} \geq 0 \\ & y_i \geq 0 \end{aligned}$$

$$\text{OPT}(\text{LP3}) \leq \text{OPT}(\text{LP}) + \lambda k$$

Called the FACILITY LOCATION PROBLEM.

Open a set S of centers and

assign points to nearest center,

but instead of asking $|S| \leq k$,
we add a cost of λ for each
open center.

$$\text{obj. fn} = \min \sum_j d(j, S) + \lambda |S|.$$

HOPE

① To find an "easy" approx algo for facility location

② Maybe for suitable λ , the algo
actually opens k centers

↓

Q Do ① & ② \Rightarrow Approx Algo is good for k -Median?

Ans This is almost true, need a little
more guarantee from the approx
algo for Facility Location.

FACILITY LOCATION OBJ. FN

↙
CONN. COST
(DISTANCE)

↘
FACILITY OPENING COST
(λ)

↗
2nd cost is what
came from the ...

Algo A is a C-LAGRANGEAN approximation
for FACILITY LOCATION if

$$\text{CONN. COST}(A) + C \cdot \text{FACILITY COST}(A) \leq C \cdot [\text{CONN}(\text{OPT}) + \text{FAC}(\text{OPT})]$$

without the C in the LHS, it is
a traditional C-approximation.

[Recall: these are minimization problems,
 $C \geq 1$].

Why does ① & ② along with C-Lagrangean
Algo \Rightarrow
good algo for k-Median?

$$\text{CONN COST}(A) + C \cdot \lambda \cdot k \leq C [\text{CONN COST}(k\text{-Median OPT}) + \lambda k]$$

Plugged in the
k-Median OPT
as a feasible solⁿ
for facility location
OPT.

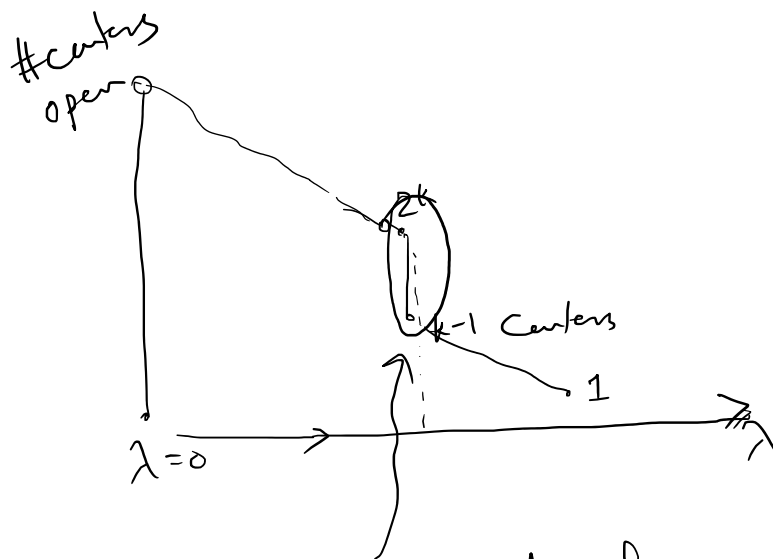


for facility location OPT.

$\text{CONN Cost}(A) \leq C \cdot \text{CONN Cost}(k\text{-Median OPT})$
 and A opens k centers (from ②)



For the suitable λ , A is a C -approx. for k -Median



What could happen in reality is that

as we keep increasing λ , there need not exist any point λ where A opens exactly k centers.

(\Rightarrow) For $\lambda - d\lambda$, it opens $> k$ centers

and $\lambda + d\lambda$ it opens $\leq k$ centers
could very well happen but there's
a very nice way to deal with it

Remains to do

- Next 1-2 lectures
- ① Design C-LAGRANGIAN algo for Facility location
 - ② Design the "Combiner procedure" for $\lambda - d\lambda$ and $\lambda + d\lambda$ problem.

26-03

Any qns about ①?

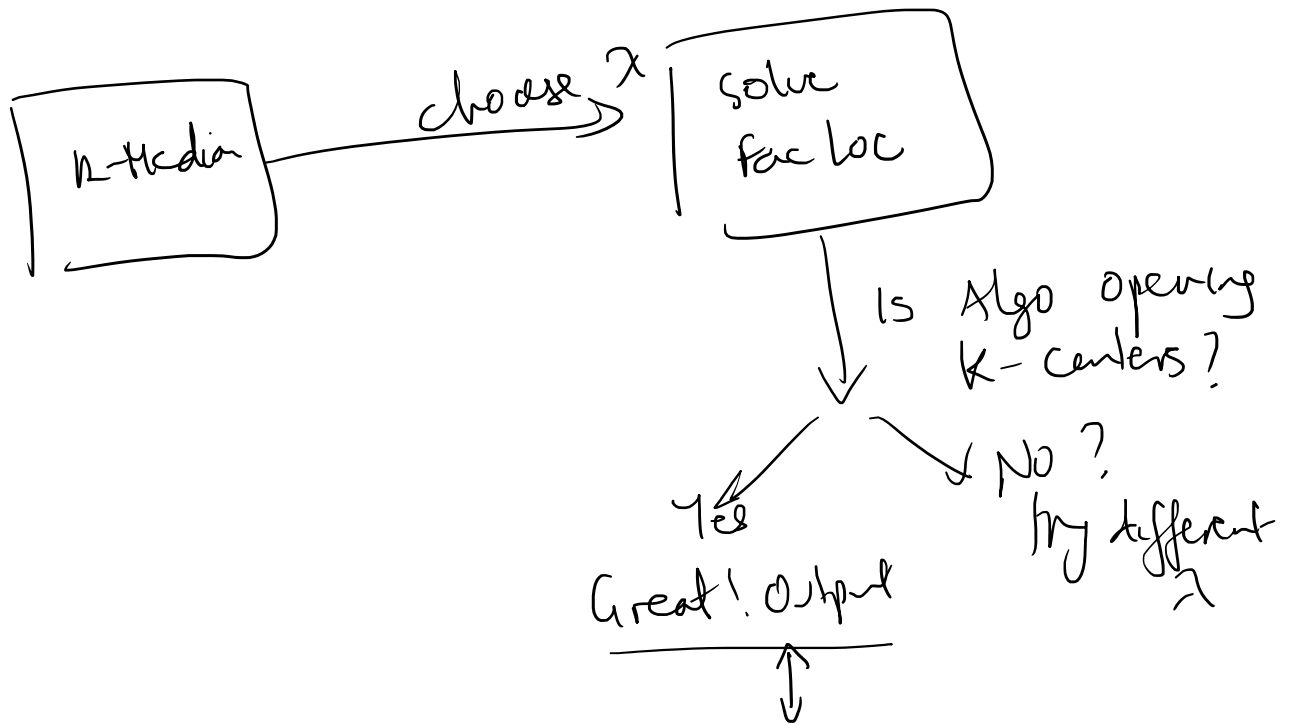
We'll assume black-box access to algorithm which has

$$\text{CONN-COST} + 3 \text{ FAC-COST} \leq 3 \text{ OPT(FAC.LOC)}$$

How do we use this for k-Median approximation?

Idea

Given a k -Median problem, let's choose an "ideal λ " and solve the Facility location instance.



Why is this solⁿ good for k -Median?

Here is where we use the

3-Lagrangian-Approx

If $\exists \lambda$ where Algo opens k centers, then

$$\begin{aligned} \text{CONN Cost (Alg)} + 3\lambda k &\leq 3(\text{OPT (FL)}) \\ &\leq 3(\text{CONN Cost (k-Median}_{\text{OPT}}) + \lambda k) \end{aligned}$$

\Rightarrow

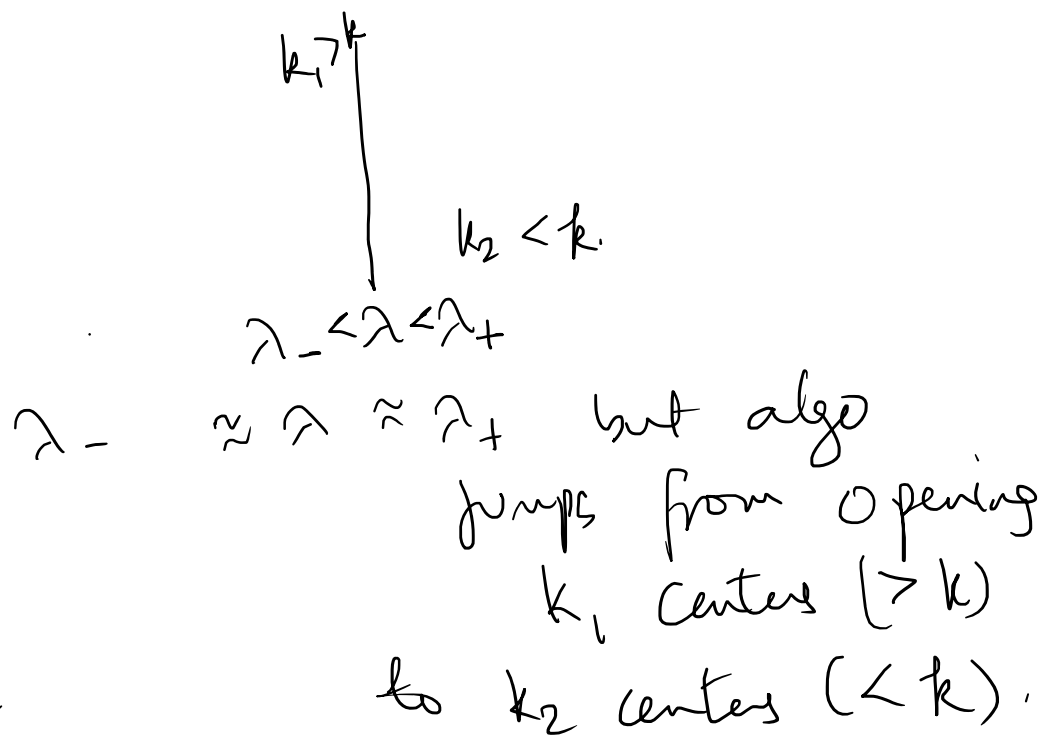
$$\text{CONN Cost (Alg)} + 3\lambda k \leq 3 \text{CONN Cost (k-Median OPT)} + 3\lambda k$$

$$\Rightarrow \text{CONN Cost (Alg)} \leq 3 \text{CONN Cost (k-Median OPT)}$$

Great, but maybe no λ is "good" where we open exactly k centers

As we keep increasing λ , Algo opens fewer centers.

↓
Maybe discrete jump occurs



Think of

$\lambda_- = \lambda_+ = \lambda$, say.
(limiting case)

$(k_1 \text{ sol}^n) \rightarrow \text{sol}^n = S_1, \text{ cost} = C_1$

$(k_2 \text{ sol}^n) \rightarrow \text{sol}^n = S_2, \text{ cost} = C_2.$

$$C_1 + 3\lambda k_1 \leq 3(\text{OPT} + \lambda k) \quad \text{--- (1)}$$

$$C_2 + 3\lambda k_2 \leq 3(\text{OPT} + \lambda k). \quad \text{--- (2)}$$

COMBINER Procedure

k_1 solⁿ is cheap but infeasible
 k_2 solⁿ is feasible but expensive
 { but, their average is feasible & Cheap }

let $p = \frac{k - k_2}{k_1 - k_2}$

consider

p (1) + $(1-p)$ (2), and see what it gives?

$$pC_1 + (1-p)C_2 + \lambda n [pk_1 + (1-p)k_2]$$

$$pC_1 + (1-p)C_2 + 3\lambda \left[\frac{pk_1 + (1-p)k_2}{\dots} \right] \leq 3(\text{OPT} + \lambda k)$$

$$pC_1 + (1-p)C_2 + 3\lambda \left[\frac{k_1k_1 - k_1k_2 + k_1k_2 - k_2k_2}{k_1 - k_2} \right] \leq 3(\text{OPT} + \lambda k)$$

$$pC_1 + (1-p)C_2 + 3\lambda k \leq 3\text{OPT} + 3\lambda k$$

So, if we choose $p = \frac{k_1 - k_2}{k_1 - k_2}$

so that

$$p \cdot k_1 + (1-p)k_2 = k \quad \left[\text{wtd Avg of } k_j \right]$$

Then corr. wtd Avg of cost

$$pC_1 + (1-p)C_2 \leq 3 \cdot \text{OPT}$$

In argument above

we used $\lambda_- = \lambda - d\lambda$ for

① and $\lambda_+ = \lambda + d\lambda$ for ②

but think of $\lambda_- = \lambda_+$ and
 $d\lambda \rightarrow 0$.

A Randomized Combiner Process :-

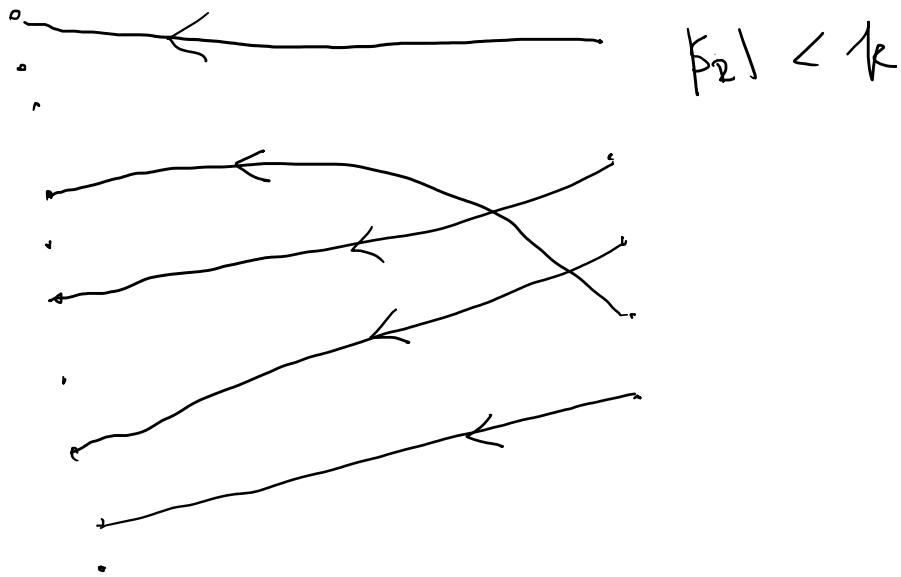
Idea :

Open centers in S_1 w.p. ' p ' and
centers in S_2 w.p. ' $(1-p)$ ' but
ensure that we open only
 k centers w.p. 1.

if we open S_1 w.p. p and
 S_2 w.p. $(1-p)$, then we are
infeasible w.p. p . Since
 $|S_1| > k$.

We 'identify good back-ups' from
 S_1 which we open instead
of all of S_1 .

$|S_1| > k$



For each $i \in S_2$, let $\eta(i) \in S_1$
denote the closest point
in S_1 to i .

$$(i.e) \quad \eta(i) = \underset{i' \in S_1}{\operatorname{argmin}} d(i, i')$$

Algo

lets assume $\eta(i_1) \neq \eta(i_2) \neq$
 $i_1, i_2 \in S_1$

Prog works even if they collide,
this is just to simplify the
discussion.

COMBINING PROCEDURE :-

- With probability $(1-p)$, Open all of S_2
 - Else, with prob p , we open $\eta(S_2)$.
- Step ①
- ↑
only open all
nodes

After step ①, we open k_2 facilities with prob 1.

From rest of S_1 ($k_1 - k_2$ points)
Choose $k - k_2$ uniformly at random and open them as centers.

⇒ After step ②, we open k centers w.p ①.

Sanity Check

Q: for $i \in S_2$, what is prob of i being selected?

↓

Ans: $(1-p)$

Q: for $i \in S_1$, what is prob of i being selected?

Ans: p

↓
Proof: if $i \in \eta(S_2)$ then it is p .

if $i \notin \eta(S_2)$ then

$$\Pr [i \text{ open}] = \frac{k_1 - k_2}{k_1 - k_2} = p.$$

LEMMA

Expected Conn. Cost of any pt j

$$\leq 2 \left[p \cdot d_1(j) + (1-p) d_2(j) \right]$$

where d_1 & d_2 are j 's cost
in S_1 and S_2 .

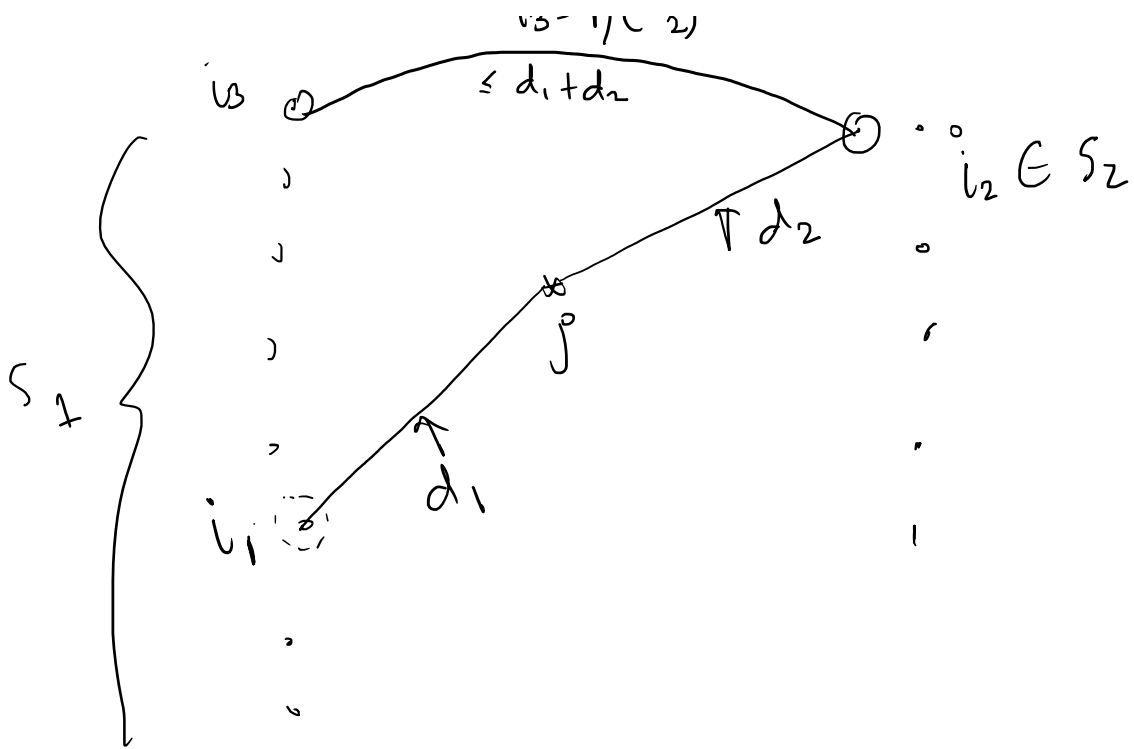
↓
Corollary

$$\text{Expected Conn. Cost of sol}^n \leq 2 \left[p \cdot C_1 + (1-p) C_2 \right]$$

$$\because \left. \begin{array}{l} C_1 = \sum_j d_1(j) \\ C_2 = \sum_j d_2(j) \end{array} \right\} \begin{array}{l} \leq 2 \cdot 3 \cdot \text{OPT} \\ = 6 \text{OPT} \end{array}$$

Proof of Lemma

$$i_3 = \eta(i_2) \\ \leq d_1 + d_2$$



Consider j and suppose its preferred conn. in S_1 is i_1 , and S_2 is i_2 .

and sps $i_3 = \eta(i_2)$ is the mate.

$$d(i_2, i_3) \leq d_1 + d_2$$

What does j connect to in our "mixed solⁿ"?

If i_1 is chosen, j can connect to it.

(happens with probability p)

If i_1 is not opened, we can check if i_2 is open.

check if i_2 is open
 if i_2 is not open, j can connect to i_3

for this "worst case analysis", let's assume that i_1 is not anybody's mate.
 $i_1 \notin N(S_2)$.

$$E[\text{Conn cost of } j]$$

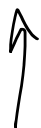
$$= p d_1 + (1-p) \cdot \left[(1-p) \cdot d_2 + p(d_1 + d_2) \right]$$

\uparrow if i_1 is open \uparrow open all of S_2 \uparrow open all of $N(S_2)$

$$= p d_1 + (1-p) [d_2 + p d_1 + p d_2]$$

$$= p d_1 + (1-p) d_2 + p(1-p)(d_1 + d_2)$$

$$\leq 2 p d_1 + 2(1-p) d_2$$



Next week

Johnson-Lindenstrauss Lemma.

Given n points X , distance function d
and opening cost of $\lambda > 0$, choose a
set S of centers to open and
connect each point to nearest open
center (incurring a cost of
 $d(j, S) = \min_{i \in S} d(i, j)$)

to minimize
$$\lambda |S| + \sum_j d(j, S)$$

PRIMAL DUAL 3-(Lagrangian) Approximation.

If opt solⁿ has facility opening cost O^*
and connection cost C^* ,
and our solution has connection cost \hat{C} ,
and facility opening cost $\hat{O} = \lambda |\hat{S}|$
if we open centers at \hat{S}

α -approximation requires that

$$\hat{C} + \hat{O} \leq \alpha (C^* + O^*)$$

α -Lagrangian Approximation (C* requires that

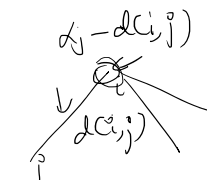


Useful to move from FL to k-Median
by choosing λ parameter
carefully.

LP (FL)	Dual (FL)
Min $\sum_j \sum_i d(i, j) x_{ij} + \lambda \sum_i z_i$	Max $\sum_j \alpha_j + \sum_{i,j} \beta_{ij}$
$\alpha_j \sum_i x_{ij} \geq 1 \quad \forall j$	$\alpha_j - \beta_{ij} \leq d(i, j)$
$\beta_{ij} \quad y_i - x_{ij} \geq 0 \quad \forall (i, j)$	$\sum_j \beta_{ij} \leq \lambda z_i$
$x_{ij}, z_i \geq 0$	$\alpha_i, \beta_{ij} \geq 0$

Recall: "take ≥ 0 linear combination, to
maximize RHS, while
ensuring all coefficients of
primal variables are dominated
by obj fn"

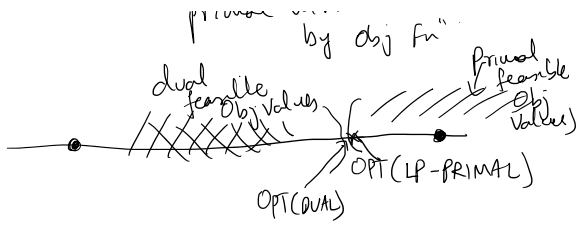
dual feasible \dots \swarrow primal feasible \dots



Think of α_j
as amount of money
 j is raising for
being connected.

β_{ij} = amt of money j
is willing to
put to opening
a center at location i

$$(\alpha_j) = \beta_{ij} + d(i, j)$$



Weak Duality THM

if (α, β) is dual feasible
and (x, y) is primal feasible,

then

$$\text{Dual Cost } (\alpha, \beta) \leq \text{Primal Cost } (x, y)$$

$$(ie) \sum_j \alpha_j \leq \sum_i \sum_j d(i, j) x_{ij} + \lambda \sum_i y_i$$

In particular,

$$(\alpha, \beta) \text{ dual feasible} \Rightarrow \boxed{\sum \alpha_j \leq OPT} \quad (*)$$

\Rightarrow if we find some good solution
with cost $\leq 3 \cdot \sum \alpha_j$, then
it will be a 3-approximation
due to $(*)$

PRIMAL-DUAL ALGORITHM : STEP 1

Initialize $\hat{T} = \emptyset$ (no open facility)
Initialize $\alpha, \beta = 0$, and all clients are
"unfrozen".

While (\exists unfrozen clients)

- Increase $\alpha_j = \alpha_j + \epsilon$ for suitably
small ϵ
for all unfrozen clients.

- If some $\alpha_j - \beta_{ij} = d(i, j)$ is
tight for some $i \in \hat{T}$, then
also increase $\beta_{ij} = \beta_{ij} + \epsilon$
to ensure ① remains feasible

- if some facility constraint ② becomes
tight, (ie), $\sum_j \beta_{ij} = \lambda$, then

\rightarrow add i to \hat{T} (ie) open i temporarily
and freeze all clients j for which
 $\alpha_j - \beta_{ij} = d(i, j)$ is tight.

Recall Dual

Max $\sum_j \alpha_j$

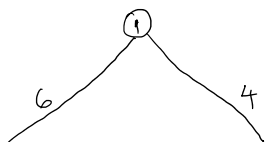
① $\alpha_j - \beta_{ij} \leq d(i, j) \quad \forall i, j$

② $\sum_j \beta_{ij} \leq \lambda \quad \forall i$

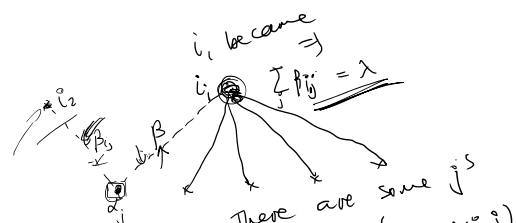
$\alpha_i, \beta_{ij} \geq 0$

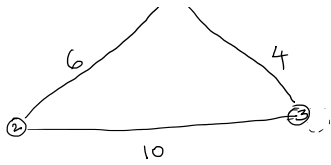
right

Toy Example
3 points, distances marked in figure



suppose $\lambda = 12$





suppon ...

There are some j s for which $d_j - \beta_{i,j} = d(i,j)$. None of these j s can $\uparrow d_j$ (all frozen)

$\hat{\tau} = \phi$, all 3 pts are unfrozen.
all increase their α slowly @ same rate

At first step itself, $\alpha_i - \beta_{ii} = d(i,i)$ is tight because $d(i,i) = 0$
 \Downarrow
 β_{ii} increases jointly with α_i

Max Z_i

$$\begin{aligned} \alpha_1 - \beta_{11} &\leq 0 \leftarrow \\ \alpha_1 - \beta_{21} &\leq 6 \leftarrow \\ \alpha_1 - \beta_{31} &\leq 4 \leftarrow \text{tight @ } t=4 \\ \alpha_2 - \beta_{12} &\leq 6 \leftarrow \\ \alpha_2 - \beta_{22} &\leq 0 \leftarrow \\ \alpha_2 - \beta_{32} &\leq 10 \leftarrow \\ \alpha_3 - \beta_{13} &\leq 4 \leftarrow \text{tight @ } t=4 \\ \alpha_3 - \beta_{23} &\leq 10 \leftarrow \\ \alpha_3 - \beta_{33} &\leq 0 \leftarrow \end{aligned}$$

Right at first step

β_{ii} grows like α_i

$$\begin{aligned} \beta_{11} + \beta_{22} + \beta_{33} &\leq 12 \leftarrow \\ \beta_{21} + \beta_{31} + \beta_{32} &\leq 12 \leftarrow \end{aligned} \quad \left. \begin{aligned} @ t=7, \beta_{11} &= 7 \\ \beta_{22} &= 1 \\ \beta_{33} &= 3 \end{aligned} \right\}$$

Think of α_j as money j is willing to raise to be connected

β_{ij} as the money j is willing to contribute to opening a facility at location i .

$$\begin{aligned} \beta_{21} &= 1 \\ \beta_{22} &= 7 \\ \beta_{23} &= 0 \end{aligned}$$

$$\begin{aligned} \beta_{31} &= 3 \\ \beta_{32} &= 0 \\ \beta_{33} &= 7 \end{aligned}$$

lets think of $\Sigma = 1$

$t=0$: All $\alpha = 0$, $\beta = 0$, $\hat{\tau} = \phi$

$t=1$: $\alpha_i = 1$, $\alpha_2 = 1$, $\alpha_3 = 1$, $\beta_{ii} = 1$, $\hat{\tau} = \phi$

$t=2$: $\alpha_i = 2$, $\beta_{ii} = 2$, $\hat{\tau} = \phi$

$t=3$: $\alpha_i = 3$, $\beta_{ii} = 3$, $\hat{\tau} = \phi$

$t=4$: $\alpha_i = 4$, $\beta_{ii} = 4$, $\hat{\tau} = \phi$

$t=5$: $\alpha_i = 5$, $\beta_{ii} = 5$, $\beta_{13} = 1$, $\beta_{31} = 1$, $\hat{\tau} = \phi$

$t=6$: $\alpha_i = 6$, $\beta_{ii} = 6$, $\beta_{13} = 2$, $\beta_{31} = 2$, $\hat{\tau} = \phi$

$t=7$: $\alpha_i = 7$, $\beta_{ii} = 7$, $\beta_{13} = 3$, $\beta_{31} = 3$, $\beta_{21} = 1$, $\beta_{12} = 1$, $\hat{\tau} = \phi$

$t=7\frac{1}{3}$: $\alpha_i = 7\frac{1}{3}$, $\beta_{ii} = 7\frac{1}{3}$, $\beta_{13} = \beta_{31} = 3\frac{1}{3}$, $\beta_{21} = \beta_{12} = 1\frac{1}{3}$

facility 1 is tight

$$\sum_j \beta_{ij} = 12$$

⇒ can't increase any β_{ij} for all j

⇒ can't increase α_j for all j st
 $\alpha_j - \beta_{ij} = d_{ij}$

⇒ freeze all such α_j .
 (In this example, all 3 clients freeze at this point)

23/03

- Think of it as a continuous process (can be discretized easily).

- Few observations

- if $\beta_{ij} > 0$ gets frozen $\alpha_j = \sum_{(i,j)} \beta_{ij} + d_{(i,j)}$ for some (i,j) pair becoming frozen/added to tight constraint
 (algs increases β_{ij} only when constraint becomes tight)

Next

if j

Then $d_{(i,j)} \leq \alpha_j$

Only those clients who can't increase their α anymore due to i freezing are frozen.

(*) these clients have

$$\alpha_j - \beta_{ij} = d_{(i,j)} \text{ is tight}$$

$$\Rightarrow \boxed{d_{(i,j)} \leq \alpha_j} \text{ since } \beta_{ij} \geq 0$$

⇒ if we open all the facilities of \hat{T} then the connection cost of all points (at the end)

$$\text{Total Conn Cost to } \hat{T} \leq \sum \alpha_j \leq \text{dual OPT} \leq \text{OPT}$$

But What about total facility cost of opening \hat{T} ?

Individually, each facility in \hat{T} is a reasonable choice to open

$$b/c \sum_j \beta_{ij} = \lambda$$

and so there are enough

clients willing to share
But issue is collective money to open it.

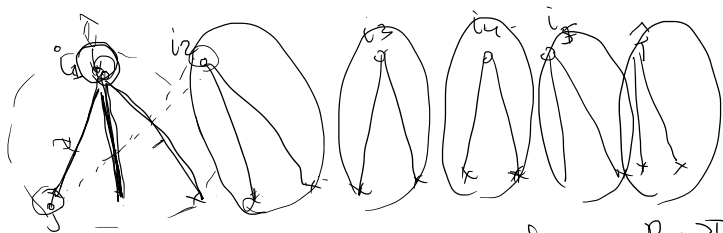
Some client could have $\beta_{ij} > 0$
to multiple facilities in \hat{T} .

GOOD CASE

For points j , there is at most 1
facility $i \in \hat{T}$ for which $\beta_{ij} > 0$

Then I claim that overall it is a
great solution

$$(10) \text{CONN Cost} + \lambda |\hat{T}| \leq \text{OPT}$$



from assumption } For each j , there is only one $\beta_{ij} > 0$
and $\text{CONN Cost}(j, \hat{T}) + \lambda |\hat{T}|$
freezes j .

$$\sum_j d_j = \sum_{\text{edges in above graph}} (d(i, j) + \beta_{ij})$$

$$= \text{CONN Cost}(j, \hat{T}) + \sum_{i \in \hat{T}} \left(\sum_{j: i} \beta_{ij} \right)$$

Hence, $\text{CONN Cost} + \text{opening Cost} \leq \sum d_j \leq \text{OPT}$

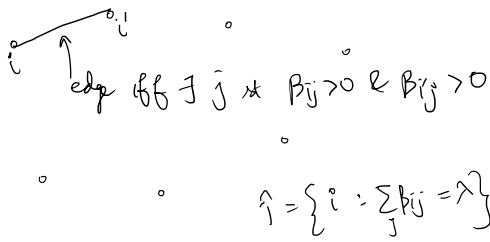
How to handle general case when some
 j is willing to put $\beta_{ij} > 0$
to multiple i 's among \hat{T}

- Form a graph with vertices in \hat{T}
- Edges (i, i') iff $\exists j$ st $\beta_{ij} > 0$
and $\beta_{i'j} > 0$

$$\& \beta_{ij} > 0$$

- Pick a maximal independent set in this graph. $T_{ind} \subseteq \hat{T}$
Means No edges amongst T_{ind}
and $\forall i \in \hat{T} \setminus T_{ind}$,
 $\exists i' \in T_{ind}$ (i, i') is edge.

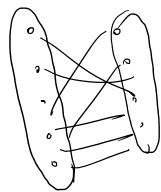
\hat{T} is all frozen centers for which $\sum_j \beta_{ij} = \lambda$



Given graph $G_1 = (V, E)$

a set of vertices $I \subseteq V$
is "INDEPENDENT SET"

iff $\forall i_1, i_2 \in I$, there
is no edge $(i_1, i_2) \in G$



example,
each side of
a bipartite
graph is an
independent set

In our case

- T_{ind} is a "maximal independent set"
- (IC) can't add any other vertex to T_{ind} while preserving independence.

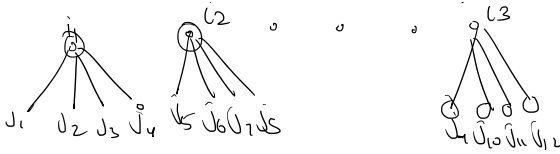
Overall Algo

- Run primal-dual process
- Build graph over \hat{T} and choose maximal independent set
- Open facilities at T_{ind} , and connect all clients to nearest open facility.

ANALYSIS

Firstly, for T_{ind} ,

the total cost of opening T_{ind} is small.



$$\lambda |T_{ind}| = \sum_{i \in T_{ind}} \left(\sum_{j: p_{ij} > 0} \beta_{ij} \right)$$

↑
j's are disjoint !!

In particular,

$$\lambda |T_{ind}| \leq \sum_j \alpha_j$$

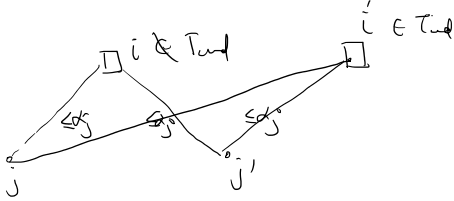
What about connection cost?



Consider client j and suppose it is frozen b/c of i

& $i \in T_{ind}$

Then $d(i, j) \leq \alpha_j$ ✓



Recap

24-03-2021

Facility location

Given points X , $|X|=n$, distance metric d , and opening cost $\lambda > 0$,

Open "centers/facilities" at $S \subseteq X$

to minimize

$$\sum_{j \in X} d(j, S) + \lambda |S|.$$

Formulated Primal-Dual and use dual to infer a solution.

$$\text{Min} \quad \quad \quad | \quad \quad \quad \text{Max} \quad \sum \lambda_i$$

$$\begin{array}{l|l}
 \text{Min} & \text{Max} \\
 \sum_j \sum_i d_{ij} x_{ij} + \sum_i \lambda y_i & \sum_j \alpha_j \\
 \sum_i x_{ij} \geq 1 \quad \forall j & \alpha_j \leq d_{ij} + \beta_{ij} + t_{ij} \\
 y_i - x_{ij} \geq 0 \quad \forall i, j & \sum_j \beta_{ij} \leq \lambda \quad \forall i \\
 x_{ij}, y_i \geq 0 &
 \end{array}$$

Intuition α_j is money j is willing to invest for its overall happiness.
 For any i , α_j breaks up into $d(i, j) + \beta_{ij}$

Share j is willing to contribute for

- ⑤ Pick a maximal independent set T_{ind} of \hat{G} .
 Alg and output that solⁿ (T_{ind})
- ① Keep raising all α_j (for unfrozen pts) until all points to nearest open center in T_{ind} .
 - ② If necessary, raise β_{ij} @ same rate.
 - ③ If facility $(\sum_j \beta_{ij} = \lambda)$ is tight for some i , freeze that i , open "temp facility" (ie) add i to \hat{T} and freeze all clients which have tight constraint $\alpha_j = \beta_{ij} + d_{ij}$.

- ④ When all clients frozen, form graph $G = (\hat{T}, \text{conflict edges})$
 (ie) (i, i') is an edge iff \exists some j with $\beta_{ij} > 0$ & $\beta_{i'j} > 0$.

For Theorem 3-Approx, we need to show:
 $\hat{C} + 3\hat{O} \leq 3(\text{OPT} + \text{FL})$ Approx.

From Lagrangian convexity, we need to show
 $\exists \lambda \geq 0$ total opening cost $\geq \lambda |T_{\text{ind}}|$
 \hat{O}

$$\hat{C} + 3\hat{O} \leq 3(\text{OPT} + \text{FL})$$

Proof

Consider T_{ind} and let us break up
 all the points in X into 2 sets:

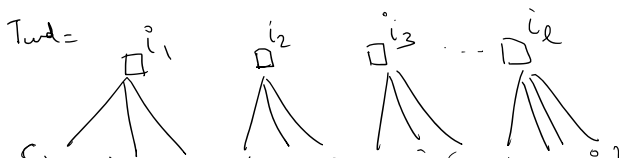
$$G = \text{good pts} = \left\{ j : \alpha_j - \beta_j = d(i, j) \text{ is tight for some } i \in T_{\text{ind}} \right\}$$

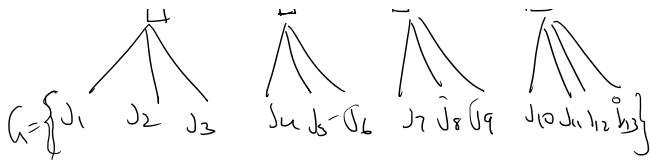
Notice: all pts j with $\beta_j > 0$ for $i \in T_{\text{ind}}$
 belong to G

Moreover,

It can't happen that $\beta_{i_1 j} > 0$
 $\& \beta_{i_2 j} > 0$
 for $i_1, i_2 \in T_{\text{ind}}$

If not, then there'll be an edge
 (i_1, i_2) so both can't be
 in independent set.





Now,

$\forall j \in G$, let $\text{mate}(j)$ be the $i \in T_{\text{ind}}$ for which $\beta_{ij} > 0$ (or if no such i exists, then pick any $i \in T_{\text{ind}}$ for which the constraint $\alpha_j - \beta_{ij} = d(i, j) = \text{tight}$)

$$\forall j \quad \alpha_j = d(j, \text{mate}(j)) + \beta(\text{mate}(j), j)$$

Sum up over all $j \in G$

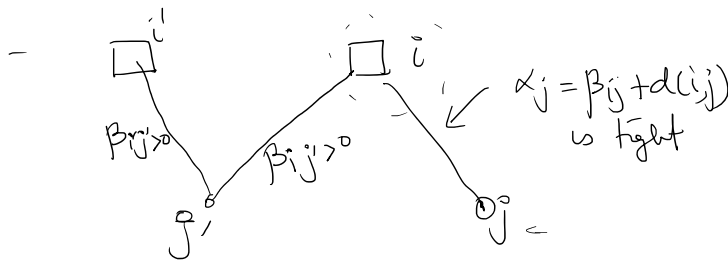
$$\begin{aligned} \sum_{j \in G} \alpha_j &= \text{CONN cost}(G) + \sum_{j \in G} \beta(\text{mate}(j), j) \\ &= \text{''} + \lambda |T_{\text{ind}}| \end{aligned}$$

$$\boxed{3 \sum_{j \in G} \alpha_j = 3 \text{CONN cost}(G) + 3 \lambda |T_{\text{ind}}|}$$

It remains to bound the bad points.

Let's look @ $j \in G$

and let the facility which "serves" j be i



Clearly $i \notin T_{\text{ind}}$ else j would have been added to T_{ind}

$\Rightarrow \exists j', i'$ st $\beta_{i'j'} > 0, \beta_{ij'} > 0$ and $i' \in T_{\text{ind}}$

Claim:

$$\textcircled{1} \quad d(i, j) \leq \alpha_j$$

$$\textcircled{2} \quad d(i, j') \leq \alpha_j$$

$$\textcircled{3} \quad d(i', j) \leq \alpha_j$$

$$\Rightarrow \boxed{d(j, i') \leq 3\alpha_j}$$

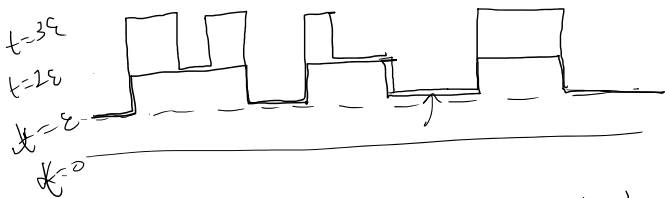
$\textcircled{1}$ is clear ($\alpha_j = \beta_{ij} + d(i, j)$
is tight).

For $\textcircled{2}$

View the dual process as
"at time t , α_j for all unfrozen
points = t "

and view $\beta_{ij} = \max(\alpha_j - d(i, j), 0)$

and if any $\sum_j \beta_{ij} = \lambda$ for some i ,
freeze all j 's for which
 $\alpha_j = \beta_{ij} + d(i, j)$.



In this view, j froze at time α_j

At this time, i gets frozen or
was already frozen.

But definitely, since at this time,

$\beta_{ij} > 0$, j' must be
frozen before j did

$$\Rightarrow \alpha_{j'} \leq \alpha_j$$

and because $\beta_{i'j} > 0$ & $\beta_{ij'} > 0$,

we get that

$$d(i, j) \leq \alpha_{j'} \leq \alpha_j$$

$$- d(i', j) \leq \alpha_{j'} \leq \alpha_j$$

$\Rightarrow j$ has a good conn to T_{ind}
of cost $\leq 3\alpha_j$

SUMMARY

$$3 \sum_{j \in G} \text{CONN Cost}(j) + 3\lambda |T_{ind}| \leq 3 \sum_{j \in G} \alpha_j$$

$$\sum_{j \in G} \text{CONN Cost}(j) \leq 3 \sum_{j \in G} \alpha_j$$

Overall

$$\text{Total CONN Cost} + 3\lambda |T_{ind}| \leq 3 \sum \alpha_j$$

$\therefore \alpha$ is feasible dual $S \leq 3(\text{OPT-FL})$

