

Complexity & Algorithms, Spring '26

Approximation:
Load Balancing
+ Clustering

Recap

- Posted next two readings

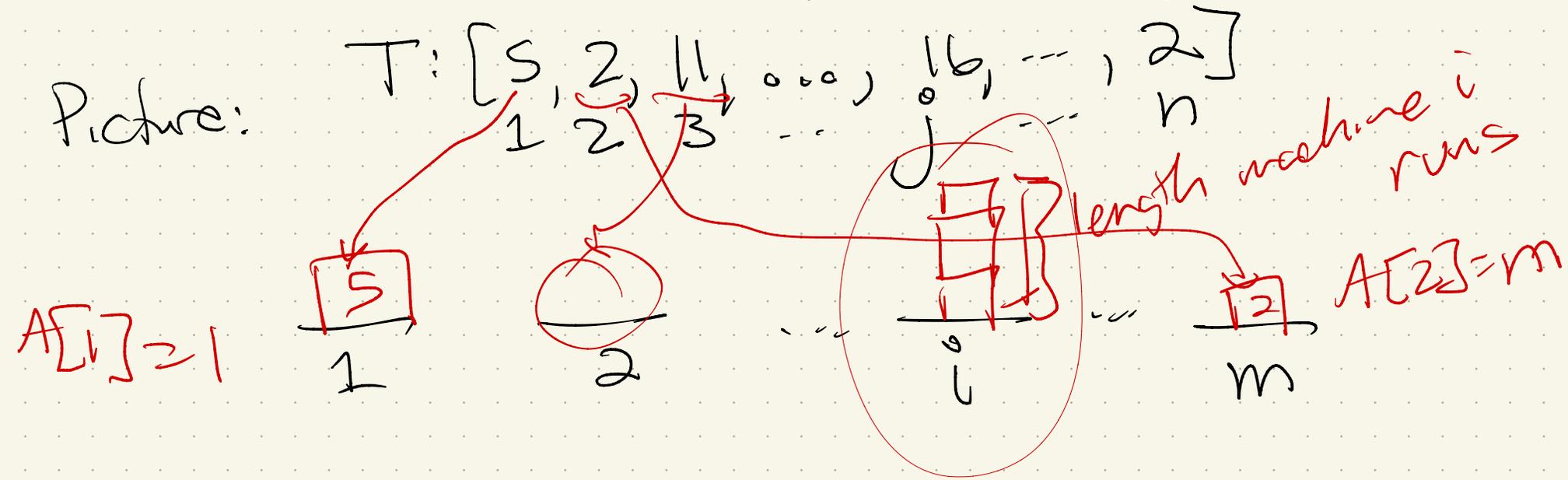
 - ↳ NP-Hardness

- HW2 - signups are on Canvas

Load Balancing

- n jobs, each with a running time $T[1..n]$
- m machines available on which to run them

Goal: Compute an assignment $A[1..n]$ where job j gets assigned to some machine $i \in [1..m]$, i.e. $A[j] = i$



Natural Goal:

Finish as early as possible!

Makespan: max time any machine is running jobs:

$$\text{makespan}(A) = \max_i \left(\sum_{j: A[j]=i} T[j] \right)$$

↑ machines
runtime

Goal: \min_A (makespan(A))

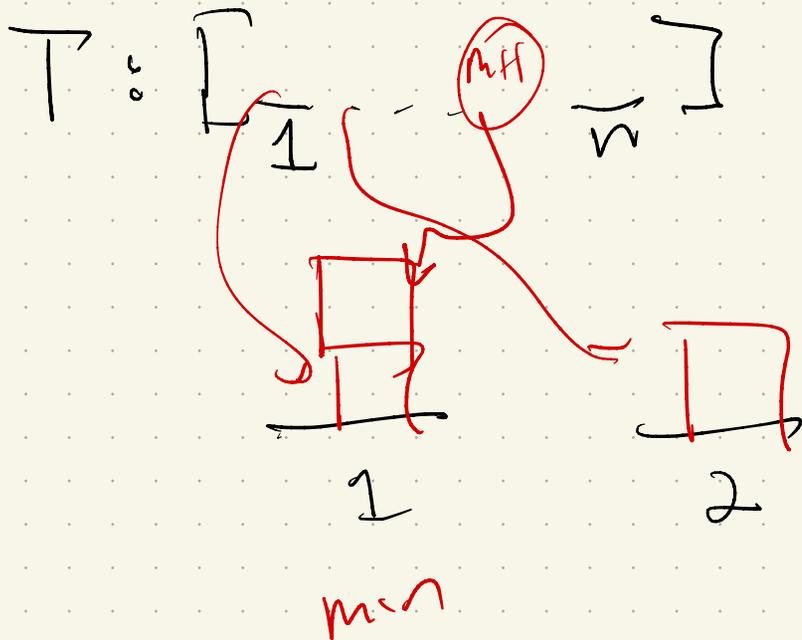
Bad news: NP-Hard

Approximating

What seems a natural strategy?

Be greedy!

Think of these
one at a time:



GREEDYLOADBALANCE($T[1..n], m$):

for $i \leftarrow 1$ to m

$Total[i] \leftarrow 0$

for $j \leftarrow 1$ to n

$mini \leftarrow \arg \min_i Total[i]$

$A[j] \leftarrow mini$

$Total[mini] \leftarrow Total[mini] + T[j]$

return $A[1..m]$

loop through jobs

} find
current
shortest
machine
& go there

Claim: Greedy make span $\leq 2 \cdot \text{OPT}$

Proof: 2 observations:

we don't know this!

① $\text{OPT} \geq \text{longest job} = \max T[j]$

if less, not running job j

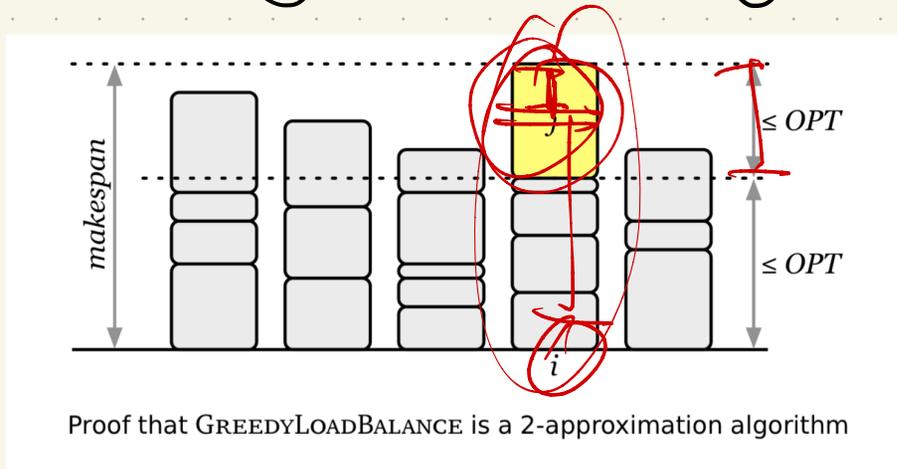
② $\text{OPT} \geq \text{average of all jobs} = \frac{1}{m} \left(\sum_j T[j] \right)$

if not:
OPT $\frac{\text{AVG}}{1} \quad \frac{1}{2} \quad \dots \quad \frac{1}{m}$

Now consider machine w / largest makespan
 in greedy alg \rightarrow machine i .

Note: $\text{Cost}(\text{greedy}) = \sum_{\text{all } A[j]=i} T[j]$

Let j be last job assigned to machine i .



Use ①: $\text{OPT} \geq \max T[j] \geq T[j]$

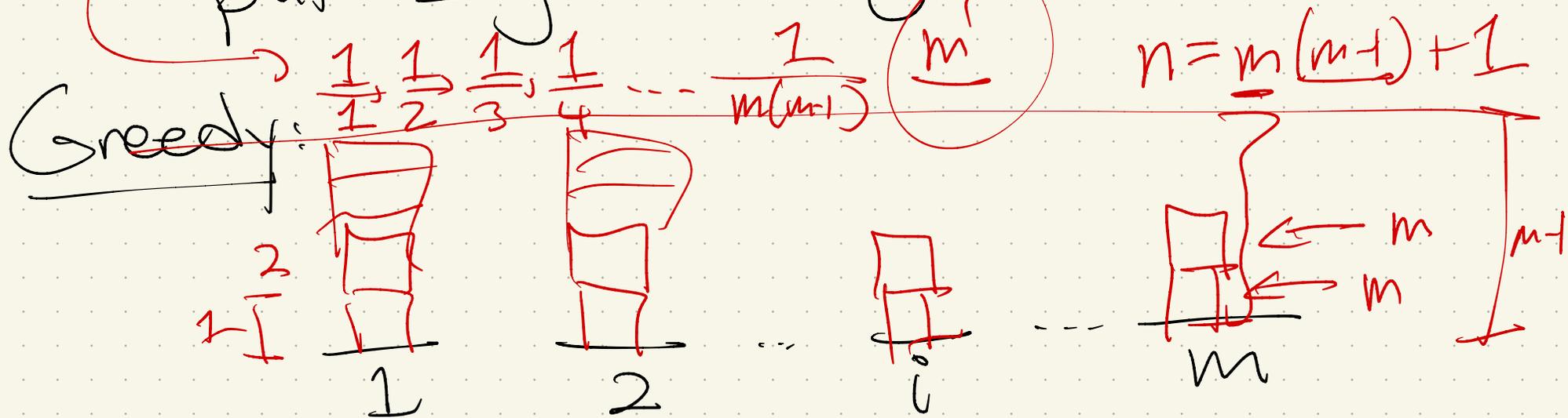
Use ②: $\text{OPT} \geq \text{avg} \geq \text{machine } i\text{'s load}$

Is this tight? adversarial!

No (but close): think lower bound

Consider $m(m-1)$ jobs of length 1

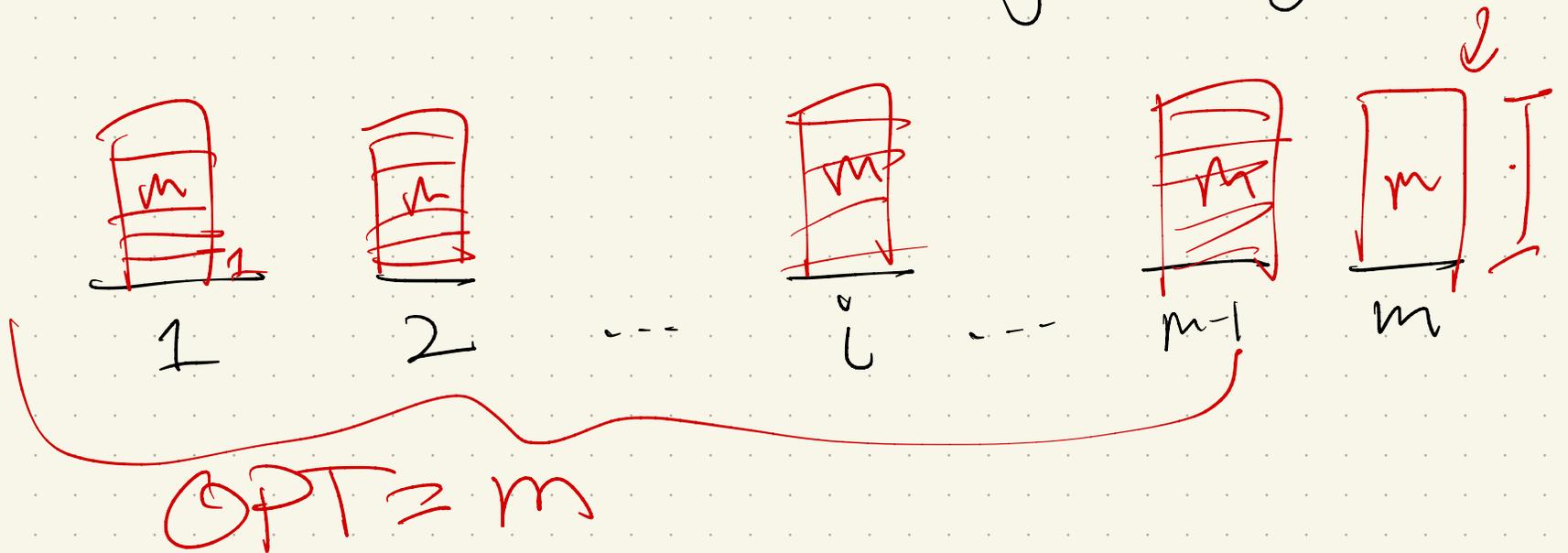
put 1 job of length m (at end)



makespan: $m-1 + m = 2m-1$

greedy

What would opt do? $m(m-1)$ jobs length 1 + 1 job length m



Ratio: $\frac{\text{Greedy}}{OPT} = \frac{2m-1}{m} < 2$

as $m \rightarrow \infty$, $\rightarrow 2$

Aside: Streaming

A model where inputs arrive one at a time.

Very common! Examples:

- processor jobs

- internet routing

⋮

Input is not known in advance

Large enough to not be able to
store everything

Improvement: Offline alg.

If we know jobs ahead of time,
can improve \rightarrow sort largest \rightarrow smallest.

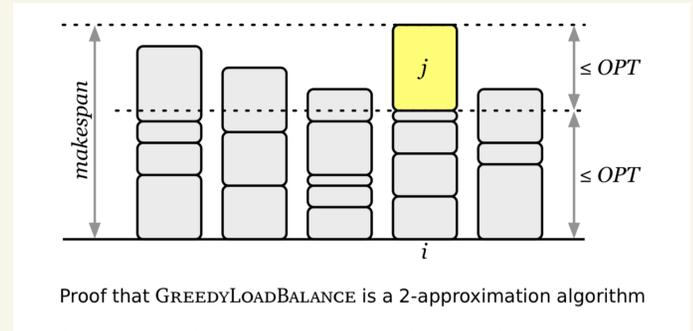
Then: • if # jobs \leq # machines:

\hookrightarrow every job goes to
new machine

• if # jobs \geq $mt + 1$:

Pigeonhole principle \Rightarrow some machine
must run 2 jobs (or more)
 $\hookrightarrow k + l$ Job

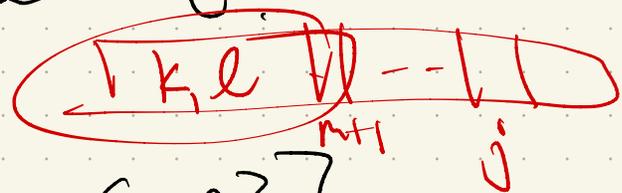
So for some $k + l \leq$ # jobs:
 $OPT \geq T[k] + T[l] \Rightarrow$ either $T[k] \leq OPT/2$
or $T[l] \leq OPT/2$



either $T[k] \geq \frac{OPT}{2}$

or $T[l] \geq \frac{OPT}{2}$

Then: Let j be last job on longest span machine i .



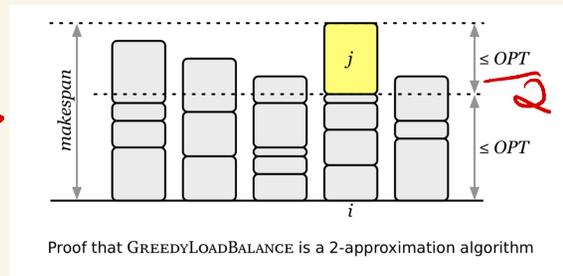
$$T[j] \leq T[m+1] \leq T[\max\{k, l\}]$$

$$= \min \{T[k], T[l]\}$$

(since we sorted big \rightarrow small)

Since $OPT \leq T[k] + T[l]$

$$\Rightarrow \min \{T[k], T[l]\} \leq \frac{OPT}{2}$$



so greedy \leq \underbrace{OPT}_{avg} + $\frac{OPT}{2}$ \uparrow last job

Question: is $3/2$ tight?

Also no! [Graham 1969] showed

$4/3$ bound. (not cover)

This one is close to tight:

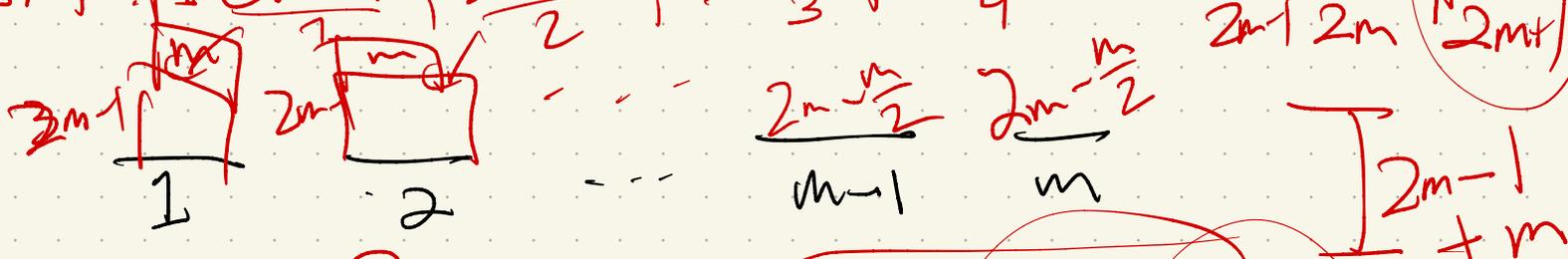
m machines, $n = 2m + 1$ jobs:

2 each of length $m, m+1, \dots, 2m-1$

plus one extra of length m

Sort: $2m-1, 2m-1, \frac{2m-2}{3}, \frac{2m-2}{4}, \dots, \frac{m}{2m-1}, \frac{m}{2m}$

greedy:

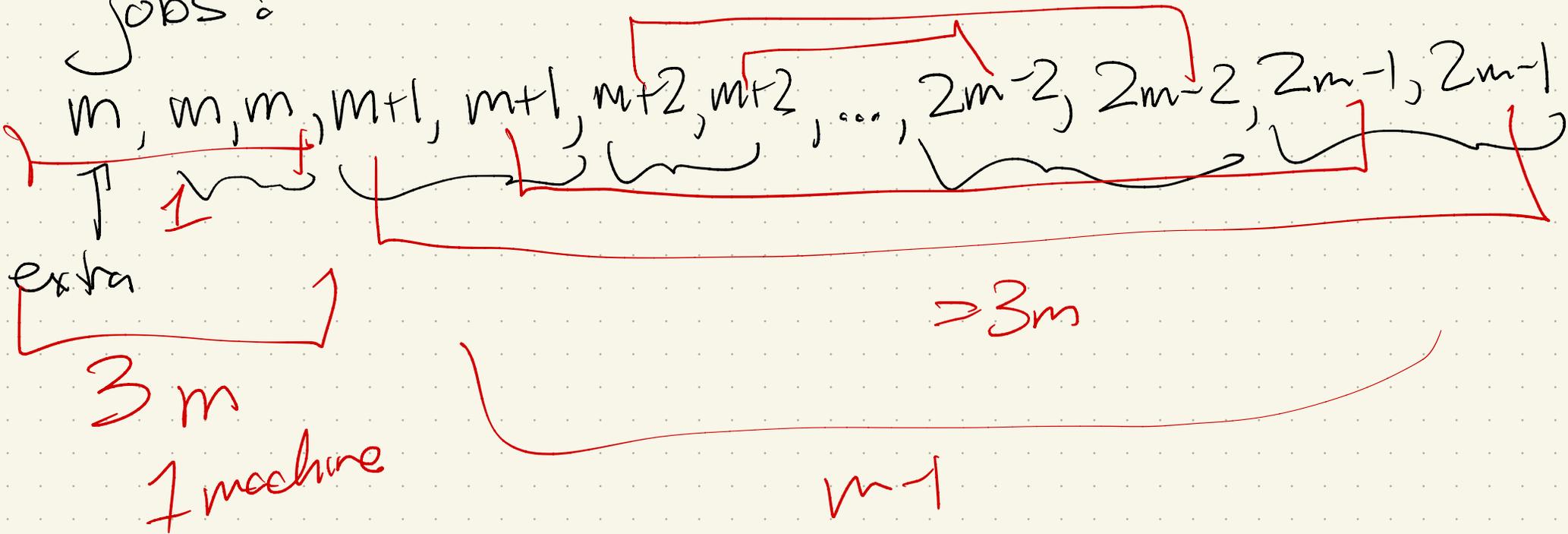


make span:

$$2m-1 + m + m = 4m-1$$

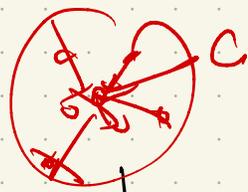
Opt: Can get $3m$:

Jobs:



So ratio: $\frac{4^{m-1}}{3^m} \approx \frac{4}{3}$

Clustering



Input: $X \subseteq \mathbb{R}^d$, $X = \{x_1, \dots, x_n\}$ + distance

$$d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

Output: A set of clusters X_1, \dots, X_k ,
each associated with a point

$$c_i \in \mathbb{R}^d: \Phi: X \rightarrow \{c_1, \dots, c_k\}$$

$$\Phi(x) = c_i$$

Cost
of cluster:

$$k\text{-means}: \min \sum_{x \in X} d(x, \Phi(x))^2$$

$$k\text{-center}: \min \left\{ \max_x d(x, \Phi(x)) \right\}$$

$$k\text{-median}: \min \sum_{x \in X} d(x, \Phi(x))$$

K-center

Bad news: NP-Hard

Even to approximate within a factor of $c < 2$

How?

Show a better than 2-approx
would let us solve some NP-Hard
problem.

K-center: Gonzalez's algorithm

Be greedy!

- Pick random point & make it first center
- While we don't have k points
choose furthest point away

Algorithm 8.2.1 Gonzalez Algorithm(X, k)

Choose $s_1 \in X$ arbitrarily. Initialize $S = \{s_1\}$.

for $i = 2$ to k **do**

Set $s_i = \arg \max_{x \in X} d(x, \phi_S(x))$.

Include $S = S \cup \{s_i\}$.

$k=3$

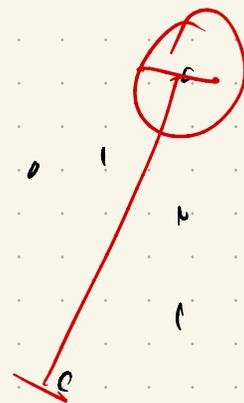
⋮

$k=4$


 s_1


 s_3


 s_2



Approximation Ratio : 2-approx

Consider end of algorithm.

Suppose we were to pick one more S_{k+1}

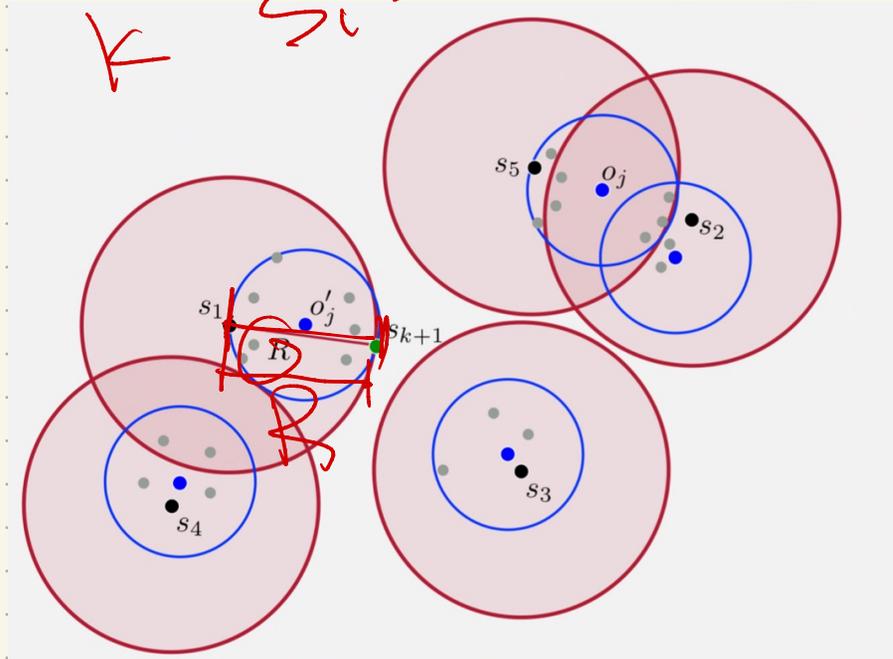
↳ some distance R away from closest center

So, $d(s_i, s_j) \geq R$ ①

Why?

If not,

I would have
chosen S_{k+1} in
my loop



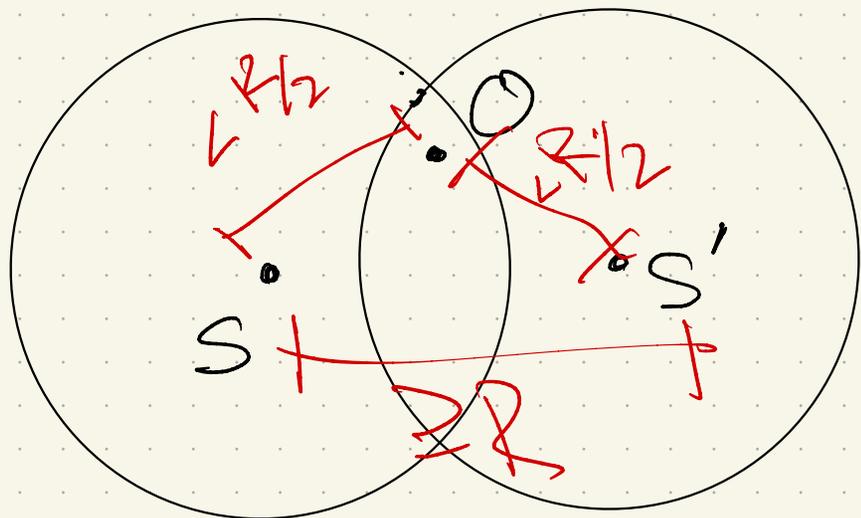
Now, proof by contradiction:

Let $\underbrace{o_1, \dots, o_k}_{\text{optimal centers}}$ be optimal centers, *

suppose their cost is $< R/2$ \rightarrow better than 2-approx

Claim: every site o must be closest to distinct s_i

Why? if not: in two R -balls



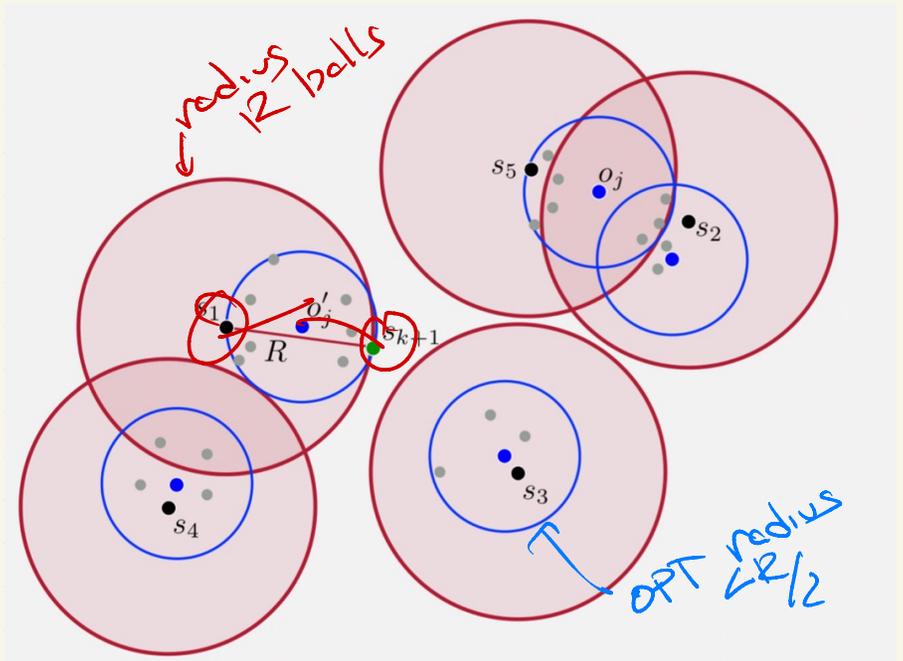
know: \triangle -ineq

$$d(s, o) + d(o, s') \geq d(s, s')$$

OK so each optimal site o_j is in a different greedy site s_j 's region!

Consider S_{k+1} again:
has $d(S_{k+1}, S_1) = R$

So, no $o_j \in \text{OPT}$
can cover both.



But all other s_i 's are $\geq R$
away from S_{k+1}

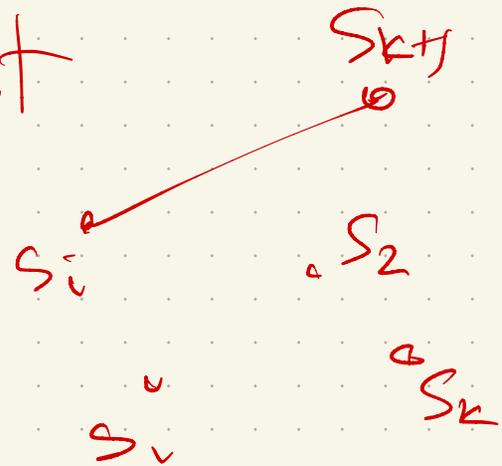
\Rightarrow $k+1$ points, no o_i can cover 2
 \hookrightarrow need $k+1$ of the o_i 's, \downarrow

\Rightarrow OPT must use $\geq R/2$ balls.

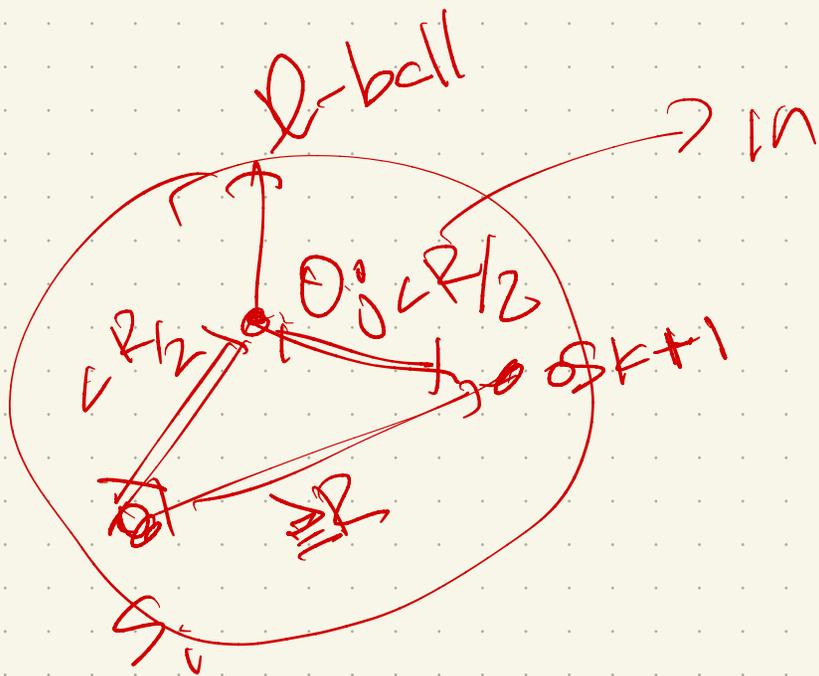
cost: \max dist { point to its center }

$\hookrightarrow R =$ greedy cost

$\hookrightarrow \text{OPT} \geq R/2$



\Rightarrow 2 approx



in an opt $< \frac{R}{2}$ solution

$$\text{OPT} = l < \frac{R}{2}$$

