# Minimum Spanning Trees / Amortized Analysis

$G = (V, E)$
$|V| = n$
$|E| = m.$

Today's lecture.

    — the Minimum Spanning Tree (MST) problem.

$$\left[ \text{Given an } \underline{\text{undirected}} \text{ graph } G = (V, E), \text{ with each edge having a weight } w_e, \text{ find a spanning tree of minimum weight.} \right]$$

(assume it for rest of lecture)

Note: if graph is connected then $\exists$ a spanning tree (a tree that contains all the vertices). Tree = no cycles, so has $n-1$ edges.

    ↳ = forest that is connected

So can enumerate over all possible spanning trees,

    for each compute the ~~cost~~ — given $T$, $w(T) = \sum_{e \in T}' w_e.$
    wt.

    and output one with least weight.

## Exponential time!

    Using the definition directly is not good here.

## A better algorithm [Kruskal '95 ]

① Sort the edges in non-decreasing order of weight.

    Say $w(e_1) \leq w(e_2) \leq \ldots \ldots \leq w(e_m).$

    Set $E_0^* \leftarrow \phi \text{ (empty set)}$

② For $i = 1$ to $m$

    if adding edge $e_i$ to $E_{i-1}^*$ does not create cycles, add it

        $(E_i^* \leftarrow E_{i-1}^* \cup \{e_i\})$

Output $T = (V, E_m^*)$

        else $E_i^* \leftarrow E_{i-1}^*$

Next Steps: ① Correctness
② Runtime.
Optional: get better algorithm, then go to step ①. :)

---

## Correctness:

Observation 1: $T_M = (V, E_{end})$ is a spanning tree  (remember: G is connected)
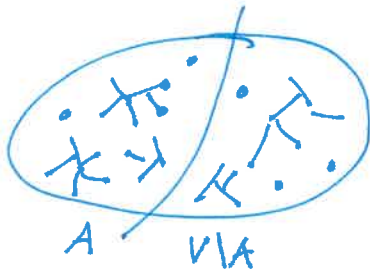(Sketch: else we would have added ~~edge~~ connecting disjoint components).

To prove $T_M$ is a min-weight spanning tree (MST), ~~let us prove useful lemma~~ Say a set of edges $S$ is __safe__

~~Lemma Pf: Suppose first~~ if $\exists$ a MST $T' = (V, E')$ such that $S \subseteq E'$.

(so we can add more edges to S to get some MST).

{
Plan: $E_0 = \phi$ is safe, trivially. Show that $E_i$ is safe $\forall i$.
}
Lets prove a general useful lemma.

Lemma 2: Suppose $S \subseteq E$ is safe. Also suppose we take a partition
$(A, V\backslash A)$ of the vertex set and $\underline{e}$ is ~~is~~ a least weight edge crossing
this, and no edge of S crosses this partition.  Then $S \cup \{e\}$ is safe.



$A$   $V\backslash A$

Pf: S is safe. So $\exists$ MST $T' = (V, E')$ st $E' \supseteq S$. If $E'$ also contains e,
done.
So $e \notin E'$. Hence adding e to spanning tree $T'$ will create a
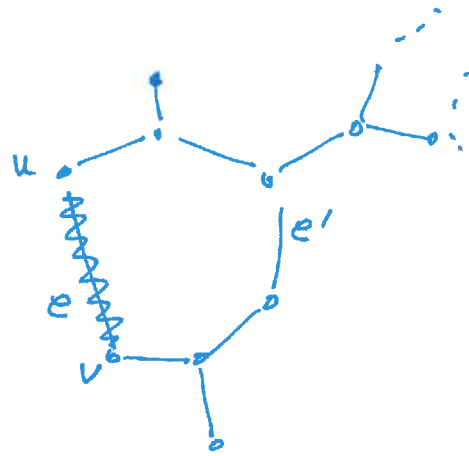cycle.

Say $e = \{u, v\}$.

Now consider the path $P$ in tree $T'$
    from $u$ to $v$. Say $u \in A$
$\exists$ an edge ~~that~~ on Path $P$     $V \in V \backslash A$.
    crosses from $A$ side to $V \backslash A$ side.
Call this edge $e'$.



Note: $w(e') \geq w(e)$ since $e$ is least-weight edge crossing
    partition $(A, V \backslash A)$

$e'$ is not in $S$    since    $S$ does not cross partition, and $e'$ does.

$T'' = \left( V, E' - \{e'\} + e \right)$ is also connected, since we
    swapped $e'$ for $e$, and both were
    on a cycle.

$\Rightarrow$ $T''$ is another spanning tree.
    and has weight at most as much as $T'$.

And it contains $S \cup \{e\}$.

     So $S \cup \{e\}$ is safe too. ☺

Good: Now using Lemma 2, we can show Kruskal produces MST.

• $E_0$ is empty set, hence trivially safe.

• if $E_{i-1}$ is safe, the edge $e_i$ does not form a cycle, so it connects two
    components of $E_{i-1}$. Call any one of them ~~$\mathbf{?}$~~ $A$, and rest is $V \backslash A$.
    Now $e_i$ is cheapest edge to cross $(A, V \backslash A)$. Using lemma 1, with $S = E_{i-1}$,
    get that $E_i = E_{i-1} \cup e_i$ is also safe. $\Rightarrow$ $E_m$ is safe by induction

$E_m$ is safe $\Rightarrow$ it can be extended into an MST.

But $E_m$ is set of edges of a spanning tree itself

$\qquad \Rightarrow$ it _is_ an MST. 🙂 .  .

_____

Aside: Can use Lemma 2 to show other algorithms. also produce MSTs.

$\qquad$ Eg. $\begin{cases} \text{Prim's Algorithm} \\[6pt] \text{Boruvka's Algo.} \end{cases}$

$\qquad$ Can discuss these some other time

_____

Good: what about run time.

Step 1. sorting edges takes $O(m \log m)$ steps.

$\qquad$ [Recall: algorithms like merge sort will do the job.

$\qquad\qquad$ quicksort ~~does~~ gets this time _in expectation_ ]

Step 2: Want that at each step, given a set

$\qquad E_i$ of edges, and a new edge $e_{i+1}$,

$\qquad\qquad$ check if $E_i \cup \{e_{i+1}\}$ contains a cycle.

$\qquad$ In other words, if $e_{i+1} = \{u, v\}$, do $u$ & $v$ lie in

$\qquad\qquad$ different connected components of $(V, E_i)$ ?

Algo 1: Just run Depth First Search on $(V, E_i)$ from $u$ to see
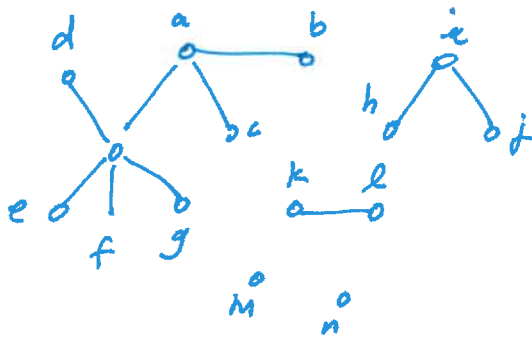
if $v$ is reachable. May take $O(i)$ time.

$\Rightarrow$ total time over all steps $\sum_{i=1}^{m} O(i) = O(m^2)$. ☹

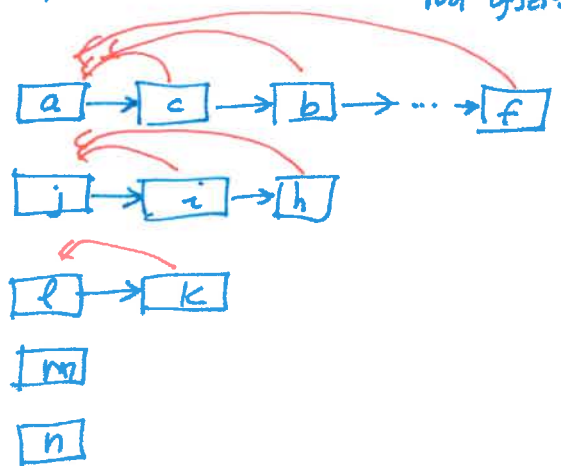Actually, always stop when $|E_i|$ has $n-1$ edges, so can

reduce to $O(mn)$.

Still bad.

---

Idea: Maintain data structures to speed up lookups.

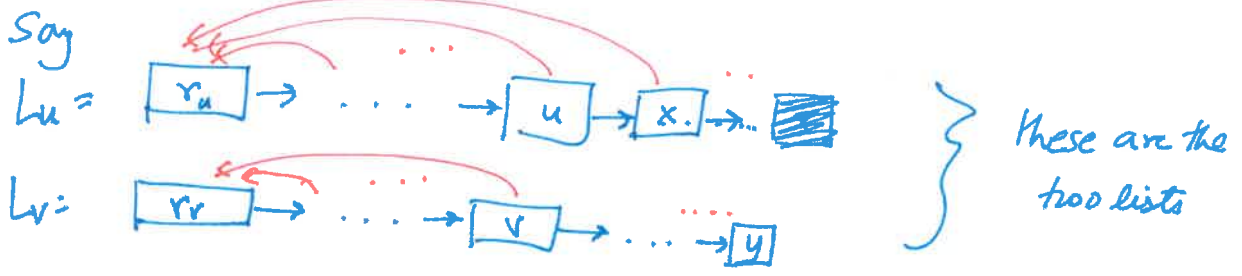~~Each step~~ Each component is a linked list, everyone maintains a pointer to "root" of set.
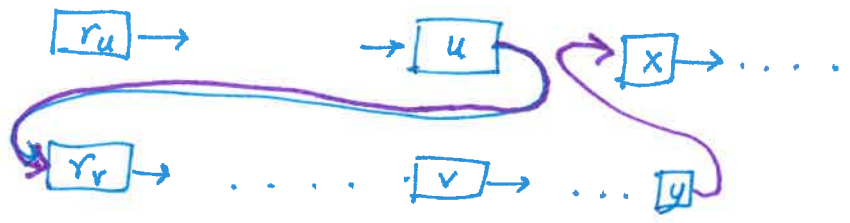


• to check if $u$ & $v$ in same component

Check: $u.root = v.root$ ?
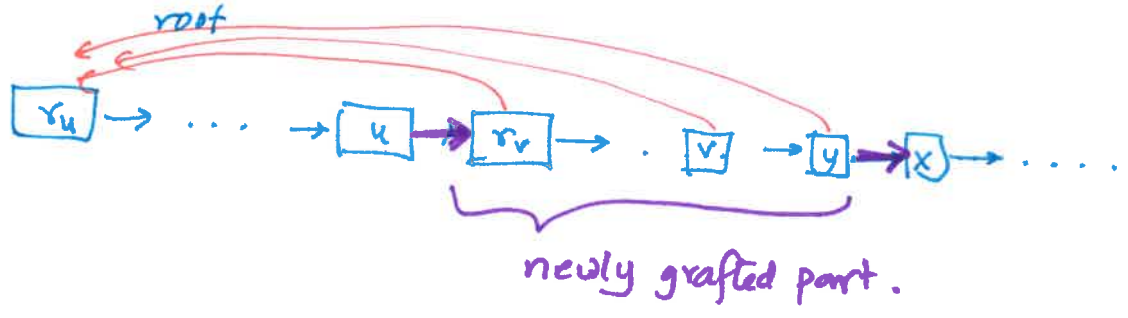
• if not, want to add $(u,v)$. So merge the components.

— need to merge the lists together.

— update root pointers.

Say

$L_u =$ 

$L_v =$

These are the two lists

Merge lists:



And then reset pointers of the second list to point to $r_u$ now.



newly grafted part.

Great. How much time?

O (length of second list).

Could be $O(n)$ in worst case.

⇒ $O(n)$ each time we add a new edge ⇒ $O(n^2)$ overall. 🙁

Final Ingredient:

When merging two lists, merge the _smaller_ into the larger.

~~Say~~ say $|L_u| \geq |L_v|$.

then "charge" each element in $L_v$ one $ each.

Use this to pay for the merge and ~~the~~ resetting pointers.

So to bound the total cost, ask:

How much does each element pay over ~~a~~the entire algo?

Each element pays 1 each time its list is merged
and it is in shorter list.

So its list length doubles each time it pays.

Final list length $\leq n$.

$\Rightarrow$ it pays $\leq \log_2 n$ times.

$\Rightarrow$ total payment $= O(n \log n)$ for merges.

$+ \; O(m)$ to check for cycles.

Kruskal's Algo : $O(m \lg m) + O(m + n \log n)$
$\underbrace{\qquad\qquad}_{\text{sorting}} \qquad \underbrace{\qquad\qquad}_{\text{checking for cycles.}}$

---

Amortized analysis: some steps cost a lot,
but on average the cost is small.

"Amortized" = "on average."

Used a ~~total~~ bank account argument
- each ~~op~~ operation was paid for by someone.
- no one paid more than $O(\log n)$.  $\Rightarrow$ total cost $\leq O(n \log n)$

\#people $\leftarrow$ max payment per person.

## Today's takeaways

- Sometimes greedy algorithms are good

  Saw it for MSTs.

  In fact this ~~extends~~ also works for a much broader class of min-weight subset problems. (See "matroids" if you are interested).

- Useful to keep data structure to avoid redoing work.

  - maintain connected components using linked lists.
  fast way to check if u & v share components.

- Choosing <u>whom</u> to merge <u>into</u> <u>whom</u> reduced time usage

  from $O(n^2)$ to $O(n \log n)$.

- Amortized analysis — clever way of accounting for the total cost.

  ( Person pays ~~the~~ only when they see improvement.

  Improves only $\leq \log_2 n$ times

  $\Rightarrow$ person pays only $\log_2 n$ times. QED )

  Will see more of it in course.