



# *Computational Discrete Optimization*

CO 353



Chaitanya Swamy

# Preface

---

**Disclaimer** Much of the information on this set of notes is transcribed directly/indirectly from the lectures of CO 353 during Winter 2022 as well as other related resources. I do not make any warranties about the completeness, reliability and accuracy of this set of notes. Use at your own risk.

Discrete optimization problems are underlying decisions that have a discrete flavor, e.g., YES/NO or  $\{0,1\}$  decisions.

The focus in this course will be on algorithms, modelling. Broad classes of problems that we will study are network connectivity problems, location problems, general integer programs.

For any questions, send me an email via <https://notes.sibeliusp.com/contact>.

You can find my notes for other courses on <https://notes.sibeliusp.com/>.

---

*Sibeliusp Peng*

# Contents

---

<b>Preface</b>	<b>1</b>
<b>1 Graph Algorithms</b>	<b>3</b>
1.1 Definitions, Notations & Terminology . . . . .	3
1.2 Shortest paths: Dijkstra's algorithm . . . . .	3
1.3 Running time and Efficient Algorithms . . . . .	5
<b>2 Graph Algorithm cont'd</b>	<b>7</b>
2.1 Minimum Spanning Trees . . . . .	7
2.2 Cut property . . . . .	8
2.3 Prim's Algorithm . . . . .	8
2.4 Kruskal's algorithm . . . . .	10
2.5 Application to Clustering . . . . .	11
<b>3 Graph Algorithms - cont'd</b>	<b>13</b>
3.1 Arborescences . . . . .	13
3.2 Min Cost Arborescence (MCA) Problem . . . . .	14
3.3 Edmond's algorithm . . . . .	15
<b>4 Matroids</b>	<b>20</b>
4.1 Introduction . . . . .	20
4.2 Max-weight independent set (MWIS) problem . . . . .	21
4.3 Applications of Matroid Optimization . . . . .	24
<b>5 Steiner Tree &amp; Computational Complexity</b>	<b>26</b>
5.1 Minimum Steiner Tree Problem . . . . .	26
5.2 MST-based algorithm . . . . .	27
5.3 Computational Complexity . . . . .	29
5.4 Polynomial-Time Reductions . . . . .	29
<b>6 P &amp; NP</b>	<b>32</b>
6.1 NP . . . . .	32
6.2 NP-hard and NP-complete . . . . .	33

# Graph Algorithms

---

## 1.1 Definitions, Notations & Terminology

A **graph** is a tuple  $(V, E)$ , where  $V$  is set of **nodes/vertices**,  $E$  is set of **edges**, where edges **joins** two nodes.

If  $e$  is an edge that joins nodes  $u, v$ , then we denote this by  $e = uv$ .  $u, v$  are called **ends** of  $e$ .  $e$  is **incident** to nodes  $u, v$ . We are not allowing parallel edges, i.e.,  $e = uv$ , and  $e' = u'v'$  are distinct edges, then  $\{u, v\} \neq \{u', v'\}$ .

An  **$u$ - $v$  path** in  $G = (V, E)$  where  $u, v \in V, u \neq v$ , is a sequence of nodes  $u_1 = u, u_2, \dots, u_k, u_{k+1} = v$ , where  $u_i u_{i+1} \in E \forall i = 1, \dots, k$ . A **cycle** in  $G$  is a sequence of nodes  $u_1, u_2, \dots, u_k, u_{k+1} = u_1$  where  $u_i u_{i+1} \in E \forall i = 1, \dots, k$ , and  $u_i$ 's are distinct. Since there are no parallel edges, we can also identify a path/cycle by its sequence of  $u_i u_{i+1}$  edges. So we will often refer to a path/cycle as a set of edges.

A graph  $G$  is **connected** if it has a  $u - v$  path  $\forall u, v \in V (u \neq v)$ .  $G$  is acyclic if  $G$  does not have a cycle. A **tree** is a connected, acyclic graph.

Let  $G = (V, E)$  be a connected graph, and  $T = (V_T, E_T)$  be a tree. IF  $E_T \subseteq E$  and  $V_T = V$ , then we say that  $T$  is a **spanning tree** of  $G$ .

If  $C$  is a cycle, and  $e \in C$ , then  $C - \{e\}$  still connects all nodes of  $C$ . So if  $G$  is a connected graph, and it contains a cycle  $C$ , and  $e \in C$ , then  $G - \{e\} := (V, E - \{e\})$  is a connected graph. Hence, a spanning tree of  $G$  is a minimal connected subgraph of  $G$ . I.e., if  $T = (V, F)$  where  $F \subseteq E$  is a minimal set such that  $(V, F)$  is connected, then  $T$  is a spanning tree of  $G$ . If  $T = (V, F)$  contains a cycle, then  $F$  is not minimal.

In **directed graph**, each edge has a direction, and goes **from** a node **to** another node.

## 1.2 Shortest paths: Dijkstra's algorithm

**Problem** Given a directed graph  $G = (V, E)$  with edge costs  $\{c_e \geq 0\}$  and a node  $s \in V$ , find the shortest path from  $s$  to all other nodes. The "shortest" path means path with the smallest total edge cost under the  $c_e$  edge costs.

**Notation** For a path  $P$ , let  $c(P) := \sum_{e \in P} c_e$  denote the total cost of  $P$ . Let  $d(u) = \min_{P: P \text{ a } s \rightarrow u \text{ path}} c(P)$ , which is shortest path (SP) distance from  $s$  to  $u$ . If  $u \rightarrow v$  is an edge of  $G$ , we have

$$d(v) \leq d(u) + c_{u,v} \quad (\clubsuit)$$

## Dijkstra's Algorithm

The idea is to maintain a set of explored vertices, and we want to expand this set. Then we can make use of  $(\clubsuit)$  to estimate the shortest path from  $s$  to  $v$ , a vertex to be added to the set. We will maintain a label  $\ell(v)$  for all  $v \notin A$ , which is our current estimate for the  $s \rightarrow v$  shortest path distance.

**Given** Directed graph  $G = (V, E)$ ,  $s \in V$ , edge costs  $\{c_e \geq 0\}$ .

---

### Algorithm 1: Dijkstra's Algorithm

---

```

1 Initialize  $A \leftarrow \{s\}$ ,  $d(s) = 0$ ,  $\ell(v) \leftarrow \infty \forall v \notin A$ .
2 while  $A \neq V$  do
3   For all  $v \notin A$  such that  $\exists u \in A$  with edge  $u \rightarrow v$ , update
      
$$\ell(v) = \min \left\{ \ell(v), \min_{u \in A: (u,v) \in E} (d(u) + c_{u,v}) \right\}$$

4   Select  $w \in V - A$  such that  $\ell(w)$  has minimum  $\ell(v)$  value among all  $v \notin A$ .
5   Update  $A \leftarrow A \cup \{w\}$ , set  $d(w) = \ell(w)$ .
```

---

#### Remark:

Can obtain actual shortest paths by maintaining along with  $\ell(w)$ , the node  $u \in A$  that determines  $\ell(w)$  (i.e.,  $u \in A$  is s.t.  $\ell(w) = d(u) + c_{u,w}$ ). Call  $u$ , the “parent” of  $w$ , and  $u \rightarrow w$  the parent edge of  $w$ .

The shortest paths obtained via previous point have a special structure: every node  $w \neq s$  has exactly one edge entering it, and there are no cycles, i.e., we have something like “directed” tree. And we denote shortest-path tree: directed tree returned by algorithm.

Also note that  $\ell(v)$  in

$$\ell(v) = \min \left\{ \ell(v), \min_{u \in A: (u,v) \in E} (d(u) + c_{u,v}) \right\}$$

is redundant, since

$$\min_{u \in A: (u,v) \in E} (d(u) + c_{u,v})$$

term only decreases as the set  $A$  only grows.

## Correctness

We may assume that there exists  $s \rightarrow u$  path in  $G \forall u \in V$ . And it's easy to modify Dijkstra's algorithm to detect if this assumption holds, and get shortest path distances from  $s$  to all nodes reachable from  $s$ .

Let  $d^{\text{Alg}}(v)$ :  $d$ -value computed by algorithm. Recall  $d(v)$  is the shortest path distance from  $s$  to  $v$ . The goal then is to show that for all  $v \in V$ ,  $d^{\text{Alg}}(v) = d(v)$ . Clearly this is satisfied when  $v = s$ .

Assume we have correctly computed shortest path distances for all  $u \in A$ ,  $\ell(v)$  is the length of the shortest path  $P$  such that *last edge of  $P$  (which enters  $v$ ) comes from a node in  $A$* .

Why? Consider such a path  $P$ . Let  $u \rightarrow v$  be the last edge of  $P$ . So  $u \in A$ ,  $d^{\text{Alg}}(u) = d(u)$ ,

$$c(P) \geq d(u) + c_{u,v} = d^{\text{Alg}}(u) + c_{u,v} \geq \ell(v)$$

and last inequality is by the definition of  $\ell(v)$ .

### Theorem 1.1

If  $w$  is added to  $A$  in line 5 of the algorithm, then  $d^{\text{Alg}}(w) = d(w)$ . (I.e., we have computed shortest path distance from  $s$  to  $w$ .)

#### Proof:

Assume we have correctly computed shortest path distance  $\forall u \in A$ . Consider an arbitrary  $s \rightarrow w$  path  $P$ . Let  $u$  be the last node on  $P$  that lies in  $A$ . Let  $v$  be the node on  $P$  after  $u$  (so  $v \notin A$ ). Let  $P'$  be the  $s \rightarrow v$  portion of  $P$ . Then

$$c(P) \geq c(P') \geq d(u) + c_{u,v} = d^{\text{Alg}}(u) + c_{u,v} \geq \ell(v) \geq \ell(w)$$

where the last equality is by the definition of  $w$  in the line 4. □

Then following parent edges gives an  $s \rightarrow w$  path of length  $= \ell(w) = d^{\text{Alg}}(w)$ .

## 1.3 Running time and Efficient Algorithms

The goal in this course is to design efficient algorithms. What does efficient mean? The short answer is “reasonable” running time.

**Running time** is number of elementary operations performed by algorithm as a function of input size. **Elementary operations** includes basic arithmetic (e.g., addition), comparisons (is  $x < y$ ?), simple logical constructs (i.e., if-then-else), assignments. **Input size** is the number of *bits* needed to specify the input. Note that number of bits need to specify a number  $x \geq 0$ ,  $x$  integer is roughly  $\log_2 x$ , which is much smaller than  $x$  itself.

For example, the size of an input of the Dijkstra’s algorithm,  $G = (V, E), \{c_e\}_{e \in E}$  is usually taken to be approximately  $|V| + |E| + \sum_{e \in E} \log_2 c_e$ .

**Reasonable running time**, i.e., efficient algorithm means that running time that is **polynomial function** of input size. In order to specify running time & input size in a convenient, compact way, we will use  $O(\cdot)$  notation.

Given two functions:  $f, g : \mathbb{R}_+ \mapsto \mathbb{R}_+$ , we say that  $f(n) = O(g(n))$  if there exist constants  $c > 0$  and  $n_0 \geq 0$  such that  $f(n) \leq c \cdot g(n)$  for all  $n \geq n_0$ .

Here are some examples:

$$\begin{aligned} n &= O(n) \\ 2n + 10 &= O(n) \\ 3n &= O(n^2) \\ \alpha n^c + \beta &= O(n^d) \\ n \log_2 n &= O(n^2) \\ \log_2 n &= O(\log_{10} n) \\ 2^n &= O(3^n) \end{aligned}$$

$f(n) = O(1)$  means  $f(n) \leq c$  for all  $n \geq n_0$ .  $f(n) = O(n^{O(1)})$  is shorthand for  $f(n)$  is bounded by some (fixed) polynomial function of  $n$ :  $f(n) \leq d \cdot n^c$ .

An algorithm with running time  $f(n)$ , where  $n$  is input size, is **efficient** if  $f(n)$  is bounded by a polynomial function of  $n$ , i.e.,  $f(n) = O(n^{O(1)})$ .

Now we can examine the running time of Dijkstra's algorithm (removing unnecessary  $\ell(v)$  in line 3). Let  $m = |E|$ ,  $n = |V|$ . We observe that there are  $n$  iterations of while loop. In each iteration:

1. Computing  $\ell(v)$  takes  $O(d^{\text{in}}(v))$  time where  $d^{\text{in}}(v)$  is the number of edges entering  $v$ .
2. Computing  $\ell(v) \forall v$  takes  $O(m)$  time since  $\sum_{v \in V} d^{\text{in}}(v) = m$ .
3. Line 4 takes  $O(n)$  time.
4. Line 5 takes  $O(1)$  time.

Each iteration takes  $O(m + n)$  time. This is  $O(m)$  if we assume there exists an  $s \rightarrow v$  path  $\forall v \in V$  since then  $m \geq n - 1$ , so  $n = O(m)$ . Then the running time of algorithm is  $O(mn)$  which is a polynomial function of input size.

However, we can have a better implementation. Observe that if  $\{u \in A : (u, v \in E)\}$  does not change across iterations, then  $\ell(v)$  does not change. So instead of recomputing  $\ell(v)$  for all  $v \notin A$ , we do the following:

When we pick  $w \notin A$  to add to  $A$ , we only update  $\ell(v)$  for all  $v \notin A$  such that  $(w, v) \in E$ , and set  $\ell^{\text{new}}(v) = \min(\ell^{\text{old}}(v), d(w) + c_{w,v})$ .

So the steps inside of while loop change as: [Let  $w^*$  be the last node added to  $A$ . Initially  $w^* = s$ .]

- (a) For every edge  $(w^*, v)$ , where  $v \notin A$ , update

$$\ell(v) = \min(\ell(v), d(w^*) + c_{w^*,v})$$

and we call this DecreaseKey operation.

- (b) Find  $w \notin A$  with minimum  $\ell(\cdot)$  value. We call this ExtractMin operation.

- (c) Update  $A \leftarrow A \cup \{w\}$ ,  $d(w) = \ell(w)$ ,  $w^* = w$ .

Across all iterations, we examine each edge  $(u, v)$  at most once in step (a) above (in the iteration when  $w^* = u, v \notin A$ ). So across all iterations,  $\leq m$  DecreaseKey operations,  $\leq n$  ExtractMin operations.

Then we can use a simple array to store  $\ell(\cdot)$  values. Note that DecreaseKey is  $O(1)$ , and ExtractMin operation is  $O(n)$ . Thus the running time =  $O(m + n^2) = O(n^2)$ .

There exist data structures such as priority queue, under which DecreaseKey and ExtractMin take  $O(\log n)$ . Then the running time is then  $O(m \log n)$ .

There exists a data structure called Fibonacci heaps, under which DecreaseKey is  $O(1)$ , and ExtractMin operation is  $O(\log n)$ . Then the running time is  $O(m + n \log n)$ .

## Graph Algorithm cont'd

---

### 2.1 Minimum Spanning Trees

**MST Problem** Given a connected, undirected graph  $G = (V, E)$ , edge costs  $\{c_e\}_{e \in E}$ . Find a spanning tree of  $G$  of minimum total edge cost.

We say “ $T$  is a spanning tree” is equivalent to “ $T$  is the edge set of a spanning tree”. We denote the cost of  $T$  by  $c(T) := \sum_{e \in T} c_e$ .

**Note:**

The  $c_e$ 's could be positive, zero, or negative.

If all  $c_e$ 's are  $\geq 0$ , then can equivalently define the MST problem as: find the min-cost connected spanning subgraph of  $G$ . Because there is always an optimal solution that is minimal connected spanning subgraph of  $G$ .

#### Theorem 2.1: Fundamental Theorem about trees

Let  $T = (V, F)$  be a graph, and  $n = |V|$ . The following are equivalent:

- (a)  $T$  is a tree (i.e., connected, acyclic)
- (b)  $T$  is connected, has  $n - 1$  edges.
- (c)  $T$  is acyclic, has  $n - 1$  edges.

**Proof:**

**(a)  $\Rightarrow$  (b)** Pick some  $r \in V$  as root node. Root  $T$  at  $r$ , i.e., draw  $T$  as hanging off of  $r$ . For each  $v \neq r$ , there is a unique edge  $uv$  of  $T$  (incident to  $v$ ) such that  $u$  is closer to  $r$  than  $v$ , and we call  $uv$  the parent edge of  $v$ .

These parent edges cover  $T$ , and number of parent edges  $= n - 1$ , since each  $v \neq r$  has a unique parent edge.

**(b)  $\Rightarrow$  (a)**  $T$  is connected. Let  $T'$  be a spanning tree of  $T$ . So by (a)  $\Rightarrow$  (b), we know that  $T'$  has  $n - 1$  edges. But  $T$  has  $n - 1$  edges. So  $T = T'$ , so  $T$  is a tree.  $\square$



## 2.2 Cut property

**Notation** Let  $v \in V$ .  $\delta(v)$  denotes the set of edges incident to  $v$ . Let  $S \subseteq V$ ,  $\delta(S) := \{uv \in E : u \in S, v \notin S\}$ . In other words,  $\delta(S)$  denotes the “boundary” of  $S$ .

**Now we assume that all edges costs are distinct.** Fix some node  $s \in V$ . Let  $e \in \delta(S)$  have the smallest edge cost among edges in  $\delta(S)$ . Is  $e$  in some MST? In fact,  $e$  is in every MST.

### cut

A cut is any partition  $(A, V - A)$  of the vertex set  $V$ , where  $A \neq \emptyset, A \subsetneq V$ .

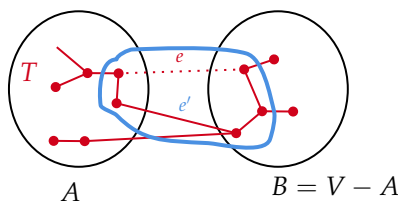
Edges crossing the cut are edges in  $\delta(A)$  ( $= \delta(V - A)$ ). If  $F \subseteq E$ , we say  $F$  crosses the cut to mean  $F \cap \delta(A) \neq \emptyset$ .

### Lemma 2.2: Cut property

Consider any cut  $(A, B)$ , where  $B = V - A$ . If  $e$  is the (unique) min-cost edge across the cut, then  $e$  belongs to every MST.

**Proof (via an exchange argument):**

Suppose  $T$  is an MST such that  $e \notin T$ . We will show that we can find another spanning tree  $T'$  (that contains  $e$ ) such that  $c(T') < c(T)$ , then a contradiction.



$T \cup \{e\}$  contains a cycle  $C$  that contains  $e$ . And this is because  $T \cup \{e\}$  is connected and has  $n$  edges, then it can't be acyclic. Here is a *basic fact*: if a cycle crosses a cut, it crosses the cut at least twice. So  $\exists e' \in C \cap \delta(A)$ ,  $e' \neq e$ . By definition of  $e$ ,  $c_e < c_{e'}$ . And  $e' \in T$ .

Consider  $T' = T \cup \{e\} \setminus \{e'\}$ . We claim that  $T'$  is a spanning tree.  $T'$  is connected since  $e' \in$  cycle in  $T \cup \{e\}$  and  $T'$  has  $n - 1$  edges. Then we have

$$c(T') = c(T) + c_e - c_{e'} < c(T)$$

a contradiction. □

## 2.3 Prim's Algorithm

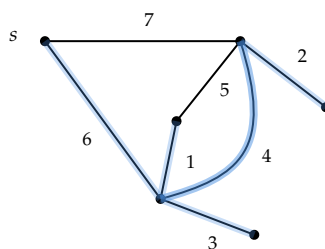
Now we can use cut property as the basis of the greedy algorithm.

### Algorithm 2: Prim's Algorithm

- 1 Pick an arbitrary “seed” node  $s \in V$ .
- 2 Initialize  $A \leftarrow \{s\}$ ,  $T \leftarrow \emptyset$ .
- 3 **while**  $A \neq V$  **do**
- 4     Choose  $e = uv \in \delta(A)$  with smallest cost, where  $u \in A, v \notin A$ .
- 5      $A \leftarrow A \cup \{v\}$ ,  $T \leftarrow T \cup \{e\}$ .
- 6 **return**  $T$

**Example:**

Blue lines are the output of Prim's algorithm.

**Theorem 2.3**

Prim's algorithm correctly computes an MST.

**Proof:**

Let  $T$  be the edge-set returned by Prim's algorithm.

- $T$  is a spanning tree.  $T$  is connected, since every node is connected to  $s$  in  $T$ . Also,  $T$  has  $n - 1$  edges. Thus  $T$  is a spanning tree.
- $T$  is a MST. Every  $e \in T$  belongs to every MST by the cut property, since it is the min-cost edge across some cut. So  $T \subseteq$  every MST. But  $T$  is itself a spanning tree, so  $T$  is MST.

□

**Corollary**

$T$  is the unique MST.

If edge costs are not distinct, then Prim's algorithm still returns an MST; there could be multiple MSTs.

**Implementation & Running Time**

Implementation will be similar to Dijkstra's algorithm. For every unexplored node  $v \notin A$ , maintain a "key"  $a(v) = \min_{e=uv: u \in A} c_e$ . So in each iteration, we choose  $w \notin A$  with smallest  $a(\cdot)$  value similar to Dijkstra, let  $w^*$  be the last node added to  $A$ . In each iteration

- For each edge  $w^*v$ , where  $v \notin A$ , update  $a(v) = \min\{a(v), c_{w^*v}\}$ . **DecKey**
- Find  $w \in V - A$  with smallest  $a(\cdot)$  value. **ExtractMin**
- Set  $A \leftarrow A \cup \{w\}$ , and  $T \leftarrow T \cup \{uw\}$ , where  $u \in A$ , and  $c_{uw} = a(w)$ .

As in Dijkstra's algorithm, across all iterations:

- $n$  **ExtractMin** operations
- $m$  **DecKey** operations

So the running time is

- $O(m + n^2)$  using a simple array to store keys (**DecKey**  $O(1)$ , **ExtractMin**  $O(n)$ )
- $O(m + n \log n)$  using a sophisticated data structure like Fibonacci Heaps (**DecKey**  $O(1)$ , **ExtractMin**  $O(\log n)$ )

## 2.4 Kruskal's algorithm

Kruskal's algorithm finds a MST. In some level, it is more greedy and intuitive than Prim's algorithm. The idea is to keep the edge costs in increasing order, and add edges to the set one by one, as long as no cycle are introduced.

---

### Algorithm 3: Kruskal's algorithm

---

```

1 Sort the edges in increasing order of cost.
2 Initialize  $T \leftarrow \emptyset$ .
3 for each edge  $e$  in sorted order do
4   if  $T \cup \{e\}$  does not have a cycle then
5      $T \leftarrow T \cup \{e\}$ 
6 return  $T$ 

```

---

### Theorem 2.4

Kruskal's algorithm returns the unique MST when all edges costs are distinct.

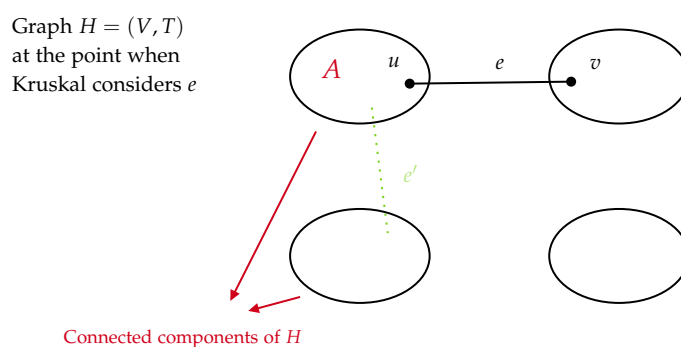
#### Proof:

Let  $T$  be the edge-set returned by Kruskal. Then,  $T$  is acyclic: by construction. Consider a basic fact: a graph is connected  $H = (V_H, E_H)$  if and only if  $\delta_H(A) \neq \emptyset \forall A : \emptyset \neq A \subsetneq V_H$ .

Suppose  $(V, T)$  is not connected. Then from the basic fact,  $\exists A, \emptyset \neq A \subsetneq V$  such that  $\delta(A) \cap T = \emptyset$ . But  $G$  is connected, so  $\exists$  some edge  $e \in \delta(A)$ . Then  $T \cup \{e\}$  is acyclic. So consider the point when Kruskal considers edge  $e$ . Let  $F \subseteq T$  be set of edges Kruskal has added until then. Then  $F \cup \{e\}$  is acyclic, so Kruskal should have added  $e$ . Then  $e \in T$ , a contradiction. So  $T$  is a spanning tree.

Consider any edge  $e = uv \in T$ . Let

$$A = \{w \in V : w \text{ is connected to } u \text{ in } T \text{ at the point when } e \text{ is considered by Kruskal}\}$$



We claim that  $e$  is the min-cost edge in  $\delta(A)$ . Observe that  $e$  is the first edge of  $\delta(A)$  considered by Kruskal. Hence  $e$  is the min-cost edge in  $\delta(A)$ . By claim,  $e \in$  every MST. So  $T \subseteq$  every MST. But  $T$  itself is a spanning tree. So  $T$  is MST.  $\square$

#### Remark:

We can stop Kruskal when  $|T| = |V| - 1$ .

**Running time** Sorting  $m$  edges takes  $O(m \log m)$ . To check if  $e = uv$  can be added, we need to check if  $u, v$  are in different components of  $(V, T)$  at that point. There exist data structures (e.g., Union-Find) for maintaining connected components that allow one to do this  $O(\log n)$  time. So total time for step 3 is  $O(m \log n)$ . Since  $m \leq n^2$ , the total time for the algorithm is  $O(m \log m + m \log n) = O(m \log n)$ .

## 2.5 Application to Clustering

**Clustering** Given a set of objects, and some notion of similarity/dissimilarity between these objects, divide the objects into groups (called clusters) so that

1. Objects in the same group are “similar” to each other.
2. Objects in different groups are “dissimilar” to each other.

**Maximum-Spanning Clustering** Given a set  $V = \{p_1, \dots, p_n\}$  of objects/points, and pairwise distances  $d(p_i, p_j) = d(p_j, p_i) \geq 0 \forall i, j \in [n]$ . The goal is to partition  $V$  into  $k$  clusters  $C_1, \dots, C_k$  ( $C_i \cap C_j = \emptyset \forall i \neq j$ ,  $\bigcup_{i=1}^k C_i = V$ ). So as to maximize the minimum inter-cluster spacing, which is equivalent to minimum distance between a pair of points in different clusters. I.e., maximize

$$\min_{\substack{i, j \in [k] \\ i \neq j}} \min_{\substack{p \in C_i \\ q \in C_j}} d(p, q)$$

---

### Algorithm 4: Single-Linkage Clustering

---

- 1 Start with every point in a separate cluster.
  - 2 Repeatedly merge the 2 clusters with smallest inter-cluster distance<sup>a</sup>, until we have  $k$  clusters.
- 

<sup>a</sup>Distance between  $C_i, C_j = \min_{p \in C_i, q \in C_j} d(p, q)$

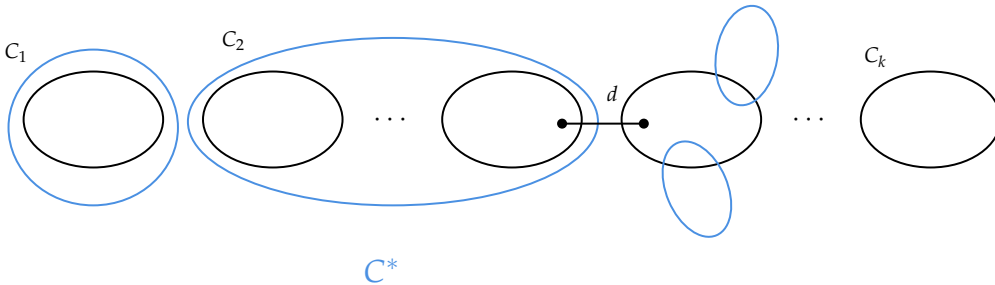
This algorithm is called Single-Linkage Clustering, which is an example of an agglomerative algorithm (i.e., based on merging clusters). Consider a graph  $G$  with  $V$  as vertex set, and an edge between every pair  $p, q \in V$ ,  $p \neq q$ , with cost  $d(p, q)$ .

**Note:**

Single-Linkage Clustering is exactly Kruskal (merging 2 clusters  $C_i, C_j$  due to points  $p \in C_i, q \in C_j$ ) when adding edge  $pq$ . *Except* that we stop when there are  $k$  components. And this is the same as taking an MST and deleting the  $k - 1$  most costly edges, i.e., the  $k - 1$  edges that Kruskal could have added last.

Thus equivalently, Run Kruskal (on complete graph with vertex set  $V$ ,  $d(p, q)$  edge costs) but stop when  $k$  components remain, which is equivalent to take MST and delete  $k - 1$  most costly edges.

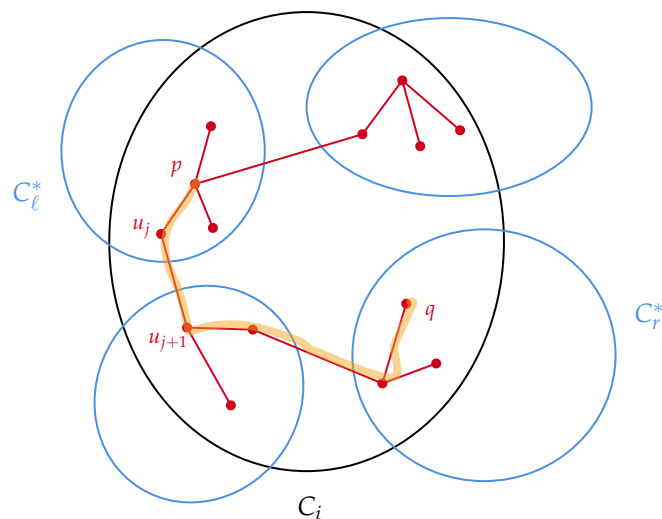
Now let's prove the correctness of Single-Linkage Clustering. Let  $C_1, \dots, C_k$  be clustering produced by MST -  $\{k - 1$  most costly edges $\}$ . Let  $d$  be the spacing of this clustering.



Observe that

$$\begin{aligned} d &= \text{cost of edge Kruskal would have added next} \\ &= (k - 1)\text{th most costly edge of MST} \end{aligned}$$

If  $C = \{C_1, \dots, C_k\}$  is not the optimum, let  $C^* = \{C_1^*, \dots, C_k^*\}$  be the optimum clustering. There might exist the case that  $C_i \subseteq C_j^*$ . As both  $C$  and  $C^*$  have  $k$  partitions, it's not possible to have  $\subseteq$  for all  $k$  partitions. Since  $C \neq C^*$ , there is some cluster  $C_i$  that intersects at least two clusters of  $C^*$ .



Red edges inside  $C_i$  all have cost  $\leq d$  since these are already added by Kruskal.

So there exists points  $p, q$  such that  $p, q \in C_i$ , but  $p, q$  lie in different clusters of  $C^*$ . Suppose  $p \in C_\ell^*$ ,  $q \in C_r^*$ ,  $\ell \neq r$ . Then considering  $p - q$  path in  $C_i$ , there must be two consecutive nodes  $u_j, u_{j+1}$  such that  $u_j \in C_\ell^*, u_{j+1} \notin C_\ell^*$ . But then spacing of  $C^* \leq d(u_j, u_{j+1}) \leq d$  since  $u_j, u_{j+1}$  are in different clusters of  $C^*$ . So this gives a contradiction since we assumed that  $\{C_1, \dots, C_k\}$  is not an optimal clustering, and optimal clustering has spacing strictly larger than  $d$ .

## Graph Algorithms - cont'd

### 3.1 Arborescences

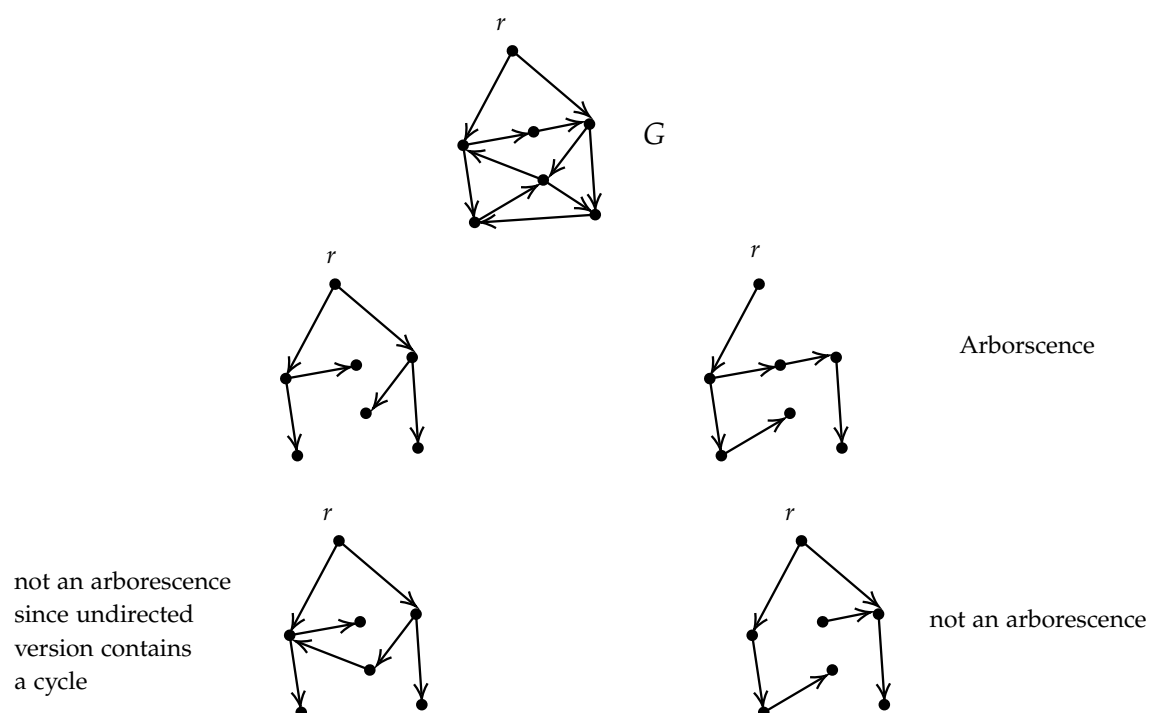
An arborescence is a directed spanning tree.  $G = (V, E)$  a directed graph, and let  $r \in V$  be a “root” node.

#### arborescence rooted at $r$

An arborescence rooted at  $r$  (or rooted out of  $r$ ) is a subgraph  $T = (V, F)$  (so  $F \subseteq E$ ) such that

- there is an  $r \rightarrow v$  path in  $T \forall v \in V$
- $T$  is a spanning tree if we ignore the directions of the edges. I.e., “undirected version of  $T$ ” is a spanning tree.

Example:

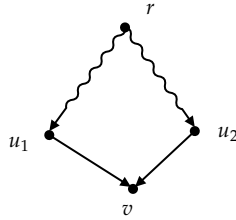


**Lemma 3.1: Useful alternate characterization of arborescence**

$T = (V, F)$  is an arborescence rooted at  $r$  if and only if  $T$  has no (directed) cycles and every  $v \neq r$  has exactly one incoming edge, (and  $r$  has no incoming edges).

**Proof:**

( $\Rightarrow$ ) Suppose  $T$  is arborescence. Since undirected version of  $T$  is a spanning tree,  $T$  has no directed cycles, and there exists  $r \rightarrow v$  path in  $T$  for all  $v \in V$  by the definition of arborescence. The path must be unique, otherwise we might encounter the situation like



then we will have a cycle in the undirected version of  $T$ . Hence every  $v \neq r$  has exactly one incoming edge.

( $\Leftarrow$ ) Suppose  $T$  has no directed cycles, and every node  $v \neq r$  has exactly one incoming edge. For any node  $v \neq r$ , we construct an  $r \rightarrow v$  path as follows: we take  $v$ 's unique incoming edge, say  $u_1 \rightarrow v$ , then  $u_1$ 's unique incoming edge, and so on... Since  $T$  has no cycle, this process must stop because we have reached  $r$ , then we found an  $r \rightarrow v$  path.

This also shows that undirected version of  $T$  is connected. Also, note that  $r$  has no incoming edges in  $T$ : suppose  $T$  had an edge  $v \rightarrow r$ , but then  $r \rightarrow v$  path in  $T$  would create a directed cycle.

So we have shown that (since every  $v \neq r$  has exactly one incoming edge)  $T$  has  $n - 1$  edges. So undirected version of  $T$  is connected, has  $n - 1$  edges, thus is a spanning tree.  $\square$

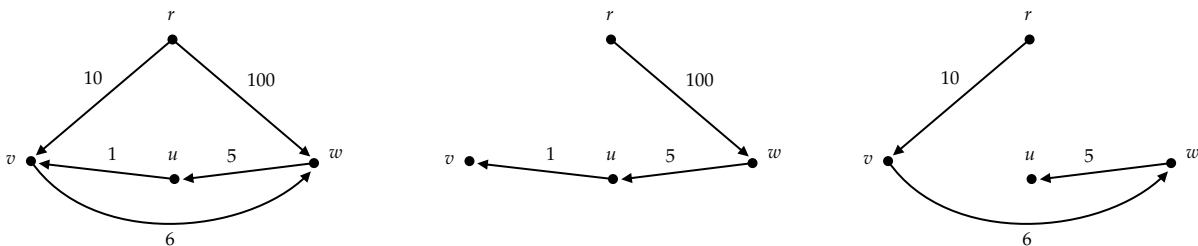
### 3.2 Min Cost Arborescence (MCA) Problem

Given a directed graph  $G = (V, E)$ , and a root node  $r \in V$ , and edge costs  $\{c_e\}_{e \in E}$ , find an arborescence rooted at  $r$  of minimum total edge cost.

**Assume there exists  $r \rightarrow v$  path in  $G$  for all  $v \in V$ , then there exists an arborescence rooted at  $r$ .**

There are two greedy strategies (inspired by MST)

**Strategy 1** Pick the cheapest edge entering a node  $v$  (for some fixed node  $v$ ). Does this edge belong to an MCA? Consider the following example for  $v$ . The cheapest edge entering  $v$  is of cost 1. However, if we pick this edge, the only arborescence containing edge  $uv$  is of cost 106, while MCA is of cost 21.



**Strategy 2** Consider some cycle, and delete the most costly edge  $e$  of the cycle. Is this valid, i.e., is there an MCA not containing  $e$ ? In the above example, we take the cycle  $C = u \rightarrow v \rightarrow w \rightarrow u$ . If we

delete the most cost edge  $vw$ , it's impossible to get MCA on the right.

Let us examine greedy strategy (1) again. Pick cheapest edge entering each node  $v \neq r$ . Let  $F^*$  be set of edges picked. Now let's make two observations.

**Observation 1** If  $(V, F^*)$  is an arborescence, then it is an MCA. Every arborescence must pick an incoming edge for every  $v \neq r$ , and  $F^*$  is the cheapest way of picking these edges. If  $(V, F^*)$  is NOT an arborescence, then by Lemma 3.1, we know that  $(V, F^*)$  contains cycle  $Z$ , not containing root  $r$ .

**Observation 2** Suppose for every node  $v \neq r$ , for each edge  $e$  entering  $v$ , we subtract a common amount  $P_v$  from its cost, i.e.,  $\forall v \neq r, \forall$  edges  $e$  entering  $v$ , we set  $c'_e = c_e - P_v$ . Then for any arborescence  $T$ , its  $c$ -cost and  $c'$ -cost differ by a constant; more precisely,  $c(T) - c'(T) = \sum_{v \neq r} P_v$  (simply because  $T$  has exactly one edge entering each  $v \neq r$  by lemma 3.1).

### Corollary

Let  $y_v = \min_{(u,v) \in E} c_{u,v}$ . For all  $v \neq r$ , for all edges  $e$  entering  $v$ , let  $c'_e = c_e - y_v$ <sup>a</sup>. Then  $T$  is an MCA with respect to  $\{c_e\}$  if and only if  $T$  is an MCA with respect to  $\{c'_e\}$  costs.

<sup>a</sup> $c'_e \geq 0$  for all  $e$ .

## 3.3 Edmond's algorithm

Now based on the corollary, imagine we have a cycle  $Z \subseteq F^*$ . And by observation 2, all edges of  $F^*$  (hence  $Z$ ) have  $c'_e = 0$ . Then as long as we can reach one node of the cycle, we can take appropriate edges of the cycle to reach all nodes in the cycle without increasing the  $c'$ -cost. As we don't care about the inside of the cycle, we can contract the cycle and run the algorithm recursively.

---

### Algorithm 5: Edmond's Algorithm for MCA

---

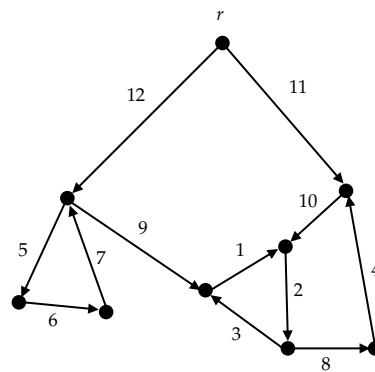
**Input:** Directed graph  $G = (V, E)$ , root  $r \in V$ ,  $\{c_e\}$  edge costs; assume there exists  $r \rightarrow v$  path for all  $v \in V$

- 0 **if**  $G$  has only one node **then return**  $\emptyset$
  - 1 **foreach**  $v \neq r$  **do**
    - $y_v := \min_{(u,v) \in E} c_{u,v}$
    - $c'_e := c_e - y_v$  for all  $e$  entering  $v$
  - 2 For each  $v \neq r$ , choose a 0  $c'$ -cost edge entering  $v$ . Let  $F^*$  be the resulting set of edges.
  - 3 **if**  $F^*$  is an arborescence **then return**  $F^*$
  - 4 Otherwise  $F^*$  contains a directed cycle. Find a cycle  $Z \subseteq F^*$  (not containing  $r$ ). Contract  $Z$  into a single supernode<sup>a</sup> to get a graph  $G' = (V', E')$ .
  - 5 Recursively find an MCA  $(V', F')$  in  $G'$  with respect to  $\{c'_e\}$  edge costs.
  - 6 Extend  $(V', F')$  to an arborescence  $(V, F)$  in  $G$ :
    - Let  $v \in Z$  be the node that has an incoming edge in  $F'$ .
    - Set  $F \leftarrow F' \cup Z \setminus \{\text{edge of } Z \text{ entering } v\}$
  - 7 **return**  $F$
- 

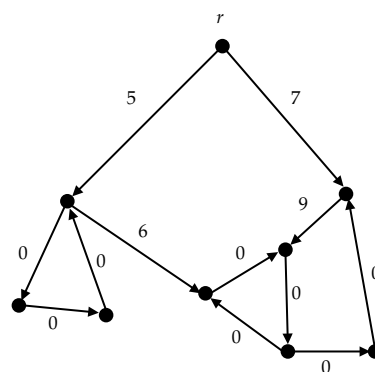
<sup>a</sup>Remove self-loops, but retain parallel edges

Now let's consider an example.

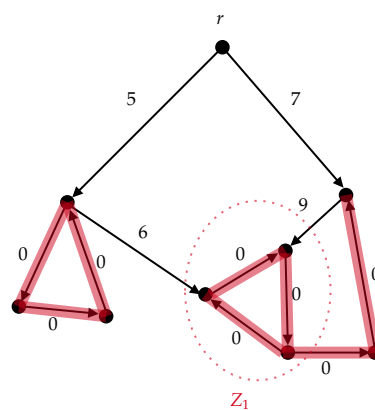




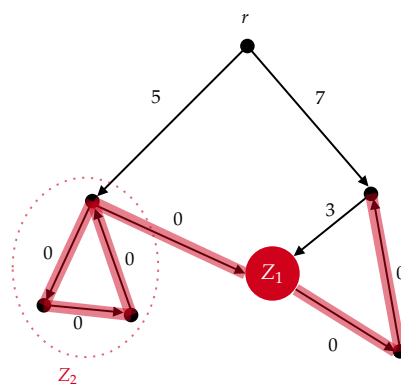
We then change the cost to  $c'_e$ .



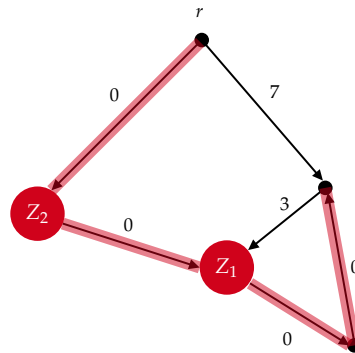
Then we mark the zero costs edges to produce  $F^*$ . And we can see there are cycles in this set of edges.



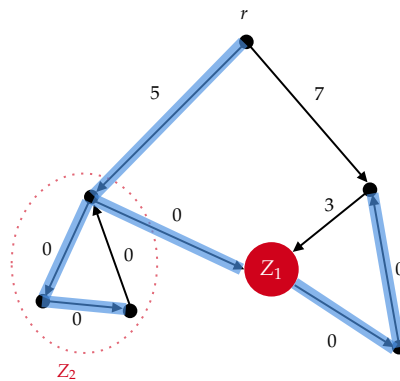
We first pick cycle  $Z_1$  and contract it. And do the cost change again.



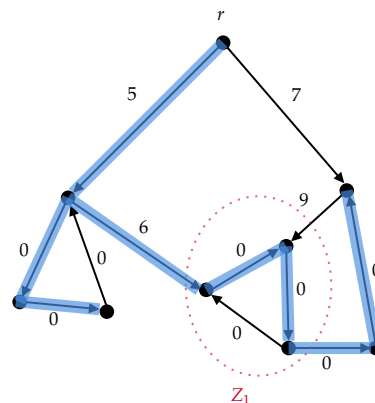
And we found cycle  $Z_2$ , and contract it, and do the cost change.



Now red edges are arborescence. Recursive calls stop. Then we expand it back.



Then expand again.



which is the output of the algorithm.

### Running time

Line 0 takes  $O(1)$ . Line 1-2 are operations on edges, take  $O(m)$ . Line 3 on checking arborescence takes  $O(n)$ : follow incoming edges in  $F^*$  to see if we can reach  $r$ , either  $rv$  path or a cycle. Line 4 takes  $O(|Z|) = O(n)$ . Line 6 takes  $O(|Z|) = O(n)$ .

Algorithm makes  $O(m)$  elementary operations and a recursive call to a smaller graph; at most  $n$  recursive calls. Thus polynomial running time.

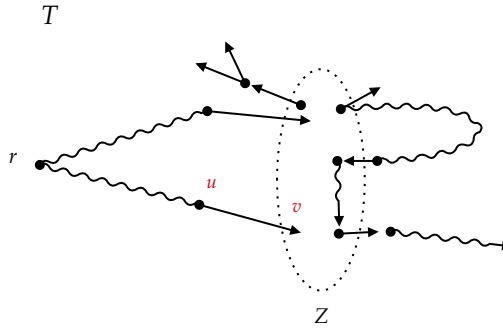
## Correctness

## Lemma 3.2

Suppose we have  $\{d_e\}$  edge costs, and a 0  $d$ -cost cycle  $Z$  such that  $r \notin Z$ . Then there exists an MCA (with respect to  $d$ -costs) that has exactly one edge entering  $Z$ .

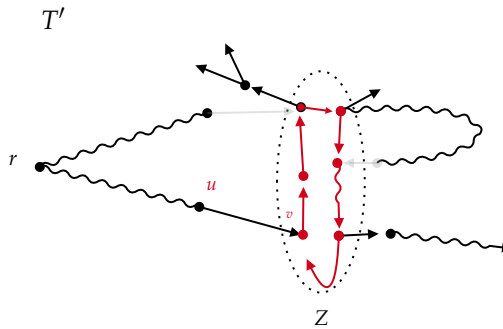
**Proof:**

Let  $T$  be any arborescence, where it might have multiple edges entering  $Z$ . We will show that  $T$  can be modified to an arborescence  $T'$  with the stated property and such that  $d(T') \leq d(T)$  where  $d(T) := \sum_{e \in T} d_e$ .  $T$  might look like the graph below.



Among all edges  $(u', v') \in T$  that enter  $Z$ , let  $(u, v)$  be such that the  $r \rightarrow u$  path in  $T$  has the fewest number of edges. Note that the  $r \rightarrow u$  path in  $T$  has no nodes from  $Z$ . Let

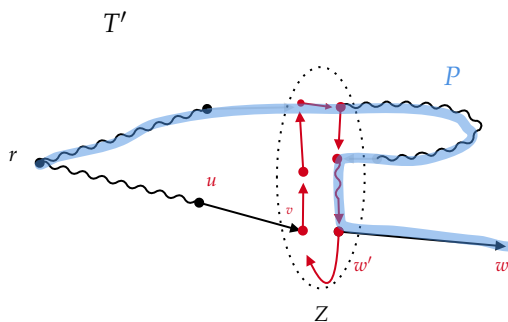
$$T' \leftarrow T - \{(u', v') \in T : v' \in Z, u' \notin Z, v' \neq v\} \cup \{\text{edges of } Z \text{ except for the edge entering } v\}$$



Note that  $d(T') \leq d(T)$  because all edges we are adding in  $Z$  have zero costs. We now have to show  $T'$  is an arborescence.

$T'$  has  $n - 1$  edges since every  $w \neq r$  has exactly one incoming edge. Now it's remained to show that there exists an  $r \rightarrow w$  path  $T' \forall w \neq r$ . Let  $P$  be the  $r \rightarrow w$  path in  $T$ . If  $P$  does not contain any node of  $Z$ , then  $P$  is also an  $r \rightarrow w$  path in  $T'$ .

Now suppose  $P$  contains a node of  $Z$ . Let  $w'$  be the last node of  $P$  in  $Z$ . Then  $T'$  contains:



- $r \rightarrow v$  path in  $T$
- $v \rightarrow w'$  portion of  $Z$
- $w' \rightarrow w$  portion of  $P$

Concatenating these gives an  $r \rightarrow w$  path in  $T'$ .

This also shows that undirected version of  $T'$  is connected, and has  $n - 1$  edges. Thus the undirected version is a spanning tree.  $\square$

**Observation** Any arborescence in  $G$  that enters cycle  $Z$  in step 4 of algorithm exactly once, yields an arborescence in  $G'$  and VICE VERSA, and these two arborescence have the same  $c'$ -cost.

### Theorem 3.3

Edmond's algorithm finds an MCA of  $G$ .

#### Proof:

By induction on  $|V|$ . Base cases where  $|V| = 1$  or  $2$  are clearly true. Suppose inductively, algorithm finds an MCA  $T'$  in  $G'$  with respect to  $c'$ -edge costs.

Let  $T^*$  be an MCA in  $G$  with respect to  $c'$ -edge costs that enters  $Z$  exactly once, which exists by lemma 3.2. Let  $T^{*'}$  be arborescence obtained from  $T^*$  for graph  $G'$ .

We have

$$\begin{aligned} c(T) &= c'(T') && \text{since } c'_e = 0 \quad \forall e \in Z \\ &\leq c'(T^{*'}) && T' \text{ is an MCA of } G' \text{ wrt. } c'\text{-edge costs, induction hypothesis} \\ &= c'(T^*) && \text{edges of } Z \text{ have } 0 \text{ } c'\text{-cost} \end{aligned}$$

So  $T$  is an MCA of  $G$  with respect to  $c'$ -edge costs. Thus  $T$  is an MCA of  $G$  with respect to  $c$ -edge costs.  $\square$

# Matroids

## 4.1 Introduction

### matroid

A matroid is a tuple  $M = (U, \mathcal{I})$ , where  $U$  is a ground set (or universe), and  $\mathcal{I} \subseteq 2^U$  is a collection of subsets of  $U$  satisfying the following properties:

- (a)  $\emptyset \in \mathcal{I}$ .
- (b) If  $A \in \mathcal{I}$ , and  $B \subseteq A$ , then  $B \in \mathcal{I}$ .
- (c) (Exchange property) if  $A, B \in \mathcal{I}$ , with  $|A| < |B|$ , then  $\exists e \in B - A$  such that  $A \cup \{e\} \in \mathcal{I}$ .

Sets in  $\mathcal{I}$  are called **independent sets**. A set not in  $\mathcal{I}$  is called a **dependent set**.

A **maximal independent set**, i.e., a set  $B \in \mathcal{I}$  such that  $B \cup \{e\} \notin \mathcal{I} \forall e \in U - B$ , is called a **basis**.

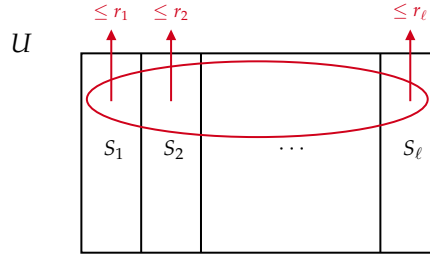
### Example: uniform matroid

Let  $U$  be any  $n$ -element set, e.g.,  $\{1, \dots, n\}$ . Let  $\mathcal{I} := \{A \subseteq U : |A| \leq k\}$ , where  $k \geq 0$ . Then  $M = (U, \mathcal{I})$  is a matroid.

- (a)  $\emptyset \in \mathcal{I}$  since  $|\emptyset| = 0 \leq k$
- (b) If  $A \in \mathcal{I}$  and  $B \subseteq A$ , then  $|B| \leq |A| \leq k$ , so  $B \in \mathcal{I}$ .
- (c) Suppose  $A, B \in \mathcal{I}$  with  $|A| < |B| (\leq k)$ . Take any  $e \in B - A$ , and note that  $|A \cup \{e\}| = |A| + 1 \leq |B| \leq k$ , so  $A \cup \{e\} \in \mathcal{I}$ .

### Example: partition matroid

Again, let  $U$  be  $n$ -element set. Let  $(S_1, \dots, S_\ell)$  be a partition of  $U$ , and let  $r_1, \dots, r_\ell \geq 0$  be some non-negative integers. Let  $\mathcal{I} := \{A \subseteq U : |A \cap S_i| \leq r_i \quad \forall i = 1, \dots, \ell\}$ .



Then  $M = (U, \mathcal{I})$  is a matroid.

(a), (b) hold trivially. For part (c), suppose  $A, B \in \mathcal{I}$  with  $|A| < |B|$ . Then there exists  $S_i$  such that  $|A \cap S_i| < |B \cap S_i|$  (since  $(S_1, \dots, S_\ell)$  partition  $U$ ). Consider any  $e \in (B \cap S_i) - (A \cap S_i)$ , and  $A' = A \cup \{e\}$ . Then

$$|A' \cap S_j| = \begin{cases} |A \cap S_j| \leq r_j & \text{if } j \neq i \\ |A \cap S_i| + 1 \leq |B \cap S_i| \leq r_i & \text{if } j = i \end{cases}$$

So  $A' \in \mathcal{I}$ .

**Example:**

Let  $U$  be collection of  $n$  vectors in  $\mathbb{R}^d$ . Let

$$\mathcal{I} := \{A \subseteq U : \text{the vectors in } A \text{ are linearly independent}\}$$

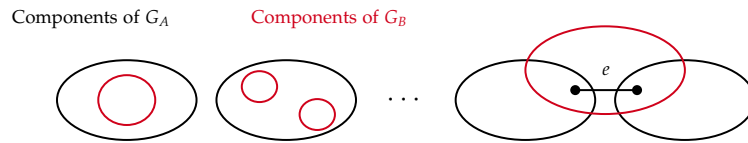
Then  $M = (U, \mathcal{I})$  is a matroid.

(a), (b) hold trivially. (c) holds because of basic linear algebra: Suppose  $A, B \in \mathcal{I}$  with  $|A| < |B|$ . If every  $v \in B$  is a linear combinations of vectors in  $A$ , then since  $|B| > |A|$ , the vectors in  $B$  must be linearly dependent. Then there exists  $v \in B$  such that  $v \notin \text{span}(A)$  (where  $\text{span}(A)$  is all vectors that are linear combinations of vectors in  $A$ ). So  $A \cup \{v\}$  consists of linearly independent vectors, i.e.,  $A \cup \{v\} \in \mathcal{I}$ .

**Example: graphic/cycle matroid**

Let  $G = (V, E)$  be an undirected graph. Let  $U = E, \mathcal{I} := \{A \subseteq E : A \text{ is acyclic}\}$ . Then  $M = (U, \mathcal{I})$  is a matroid.

(a), (b) hold trivially. Consider (c). Suppose  $A, B \in \mathcal{I}$  with  $|A| < |B|$ . Then  $G_A = (V, A)$  has  $n - |A|$  components where  $n = |V|$ , and  $G_B = (V, B)$  has  $n - |B| < n - |A|$  components.



In the picture above, we denote the components of  $G_B$  with red color. Some components of  $G_B$  can be contained within components of  $G_B$ . However, it cannot be that the vertex-set of every component of  $G_B$  is a subset of the vertex-set of some components of  $G_A$ . I.e., there exists some component of  $G_B$  that intersects at least two components of  $G_A$  as shown in the graph above.

This means there exists some  $e \in B$  that connects two components of  $G_A$ . Thus  $A \cup \{e\}$  is acyclic, so  $A \cup \{e\} \in \mathcal{I}$ .

## 4.2 Max-weight independent set (MWIS) problem

Here we present two matroid optimization problems.

**Max-weight independent set (MWIS) problem**

Given a matroid  $M = (U, \mathcal{I})$ , and weights  $\{w_e\}_{e \in U}$  (the  $w_e$ 's could be arbitrary), find a max-weight independent set, i.e., find  $A \in \mathcal{I}$  such that

$$w(A) := \sum_{e \in A} w_e = \max_{B \in \mathcal{I}} w(B)$$

**Matroid intersection problem**

Given two matroids  $M_1 = (U, \mathcal{I}_1)$  and  $M_2 = (U, \mathcal{I}_2)$  and weights  $\{w_e\}_{e \in U}$ , find a max-weight set that is independent in both matroids (common independent set), i.e., solve

$$\max_{A \in \mathcal{I}_1 \cap \mathcal{I}_2} w(A)$$

Matroid intersection problem is beyond the scope of this course, but you can learn it from [CO 450](#).

Note that *max-weight spanning tree* with positive edge weights is a special case of MWIS problem, where the matroid is the graphic matroid. (since bases of graphic matroid associated with a connected graph are spanning trees). MST with  $\{c_e\}$  edge costs can be captured by max-weight spanning tree with positive edge weights. By defining  $w_e = M - c_e$ , where  $M > \max_e c_e$  (so that  $w_e > 0$ ). For any spanning tree  $T$   $w(T) = (n-1)M - c(T)$ .

Now we develop the greedy algorithm for MWIS. Input is  $M = (U, \mathcal{I})$ ,  $\{w_e\}_{e \in U}$ . We may assume that  $w_e > 0$  for all  $e \in U$ . Because otherwise, we can move to the smaller matroid

$$M' = (U' := \{e \in U : w_e > 0\}, \mathcal{I}' = \{A \subseteq U' : A \in \mathcal{I}\})$$

where we can verify  $M'$  is a matroid. Solving MWIS on  $M'$  will also solve MWIS on  $M$ . Since every independent set of  $M'$  is also independent in  $M$ . If  $A \in \mathcal{I}$ , then  $A \cap U' \in \mathcal{I}'$  and  $w(A \cap U') \geq w(A)$ .

**Algorithm 6:** Greedy Algorithm for MWIS

**Input:**  $M = (U, \mathcal{I})$ ,  $\{w_e\}_{e \in U}$

- 1 Sort elements in decreasing order of weight.
- 2 Initialize  $A \leftarrow \emptyset$
- 3 Considering elements in sorted order, if  $A \cup \{e\}$  is independent, then set  $A \leftarrow A \cup \{e\}$  where  $e$  is current element being considered.
- 4 **return**  $A$

Observe that above algorithm run on graphic matroid is equivalent to Kruskal for max-weight spanning tree.

**Correctness****Claim 4.1**

If  $B, B'$  are two bases of a matroid, then  $|B| = |B'|$ .

**Proof:**

Suppose not, and  $|B| < |B'|$ . Then since  $B, B' \in \mathcal{I}$ , by the exchange property of matroids, there exists  $e \in B' - B$  such that  $B \cup \{e\} \in \mathcal{I}$ , contradicting that  $B$  is a maximal independent set.  $\square$

Now let  $A$  be the set returned by the greedy algorithm.

**Claim 4.2**

$A$  is a basis of  $M$ .

**Proof:**

Suppose there exists  $e \notin A$  such that  $A \cup \{e\} \in \mathcal{I}$ . Then consider the point when  $e$  is considered by greedy. At that point, we have some set  $S \subseteq A$ . But then  $S \cup \{e\} \in \mathcal{I}$  since  $S \cup \{e\} \subseteq A \cup \{e\}$ , so algorithm should have added  $e$ , contradicting that  $e \notin A$ .  $\square$

**Theorem 4.3**

Greedy algorithm returns a max-weight independent set.

**Proof:**

Let  $A^*$  be a max-weight independent set. We want to show that  $w(A) = w(A^*)$ , hence  $A$  is a max-weight independent set.

Observe that  $A^*$  is a basis of  $M$ . Otherwise if there exists  $e \notin A^*$  such that  $A^* \cup \{e\} \in \mathcal{I}$ , we have  $w(A^* \cup \{e\}) > w(A^*)$ . So we have  $|A| = |A^*|$  by Claim 4.1 and Claim 4.2.

Let  $k = |A| = |A^*|$ . Suppose  $w(A) < w(A^*)$ . Let

$$\begin{aligned} A &= \{e_1, e_2, \dots, e_k\} \\ A^* &= \{e_1^*, e_2^*, \dots, e_k^*\} \end{aligned}$$

where elements are ordered by the ordering used by greedy.

Let  $A_i := \{e_1, \dots, e_i\}$  and  $A_i^* := \{e_1^*, \dots, e_i^*\}$  for all  $i = 1, \dots, k$ . And we define  $A_0 = A_0^* = \emptyset$ . Consider the smallest index  $j$  such that  $w(A_j) < w(A_j^*)$ . Such  $j$  exists since  $w(A_k) = w(A) < w(A^*) = w(A_k^*)$ . We have

- $|A_j^*| = j$ ,  $A_j^* \in \mathcal{I}$  since  $A_j^* \subseteq A^*$
- $|A_{j-1}| = j-1$ ,  $A_{j-1} \in \mathcal{I}$  since  $A_{j-1} \subseteq A$
- By the exchange property, there exists  $e \in A_j^* - A_{j-1}$  such that  $A_{j-1} \cup \{e\} \in \mathcal{I}$ .

Since  $e \in A_j^*$  and  $e_j^*$  has the least weight among elements of  $A_j^*$ , we have  $w_e \geq w_{e_j^*}$ . Since

$$w(A_{j-1}) + w_{e_j} = w(A_j) < w(A_j^*) = w(A_{j-1}^*) + w_{e_j^*}$$

and  $w(A_{j-1}^*) \leq w(A_{j-1})$ . Thus  $w_{e_j^*} > w(e_j)$ . To summarize, we have

$$w_e \geq w_{e_j^*} > w_{e_j}$$

Consider the point when greedy considers element  $e$ . At this point, we have some set  $S \subseteq A_{j-1}$ . This is because  $w_e > w_{e_j}$ , so greedy must consider  $e$  before  $e_j$ . But then  $S \cup \{e\} \in \mathcal{I}$  since  $S \cup \{e\} \subseteq A_{j-1} \cup \{e\}$  and so greedy should have added  $e$ . Contradiction, since then we would have  $e \in A_{j-1}$ .  $\square$

**Running Time and Input Specification**

How is the matroid  $M$  given as input? It would be space inefficient if we give all sets in  $\mathcal{I}$  as input. So  $M$  is specified by means of the universe  $U$  and a **matroid independence oracle**, which is a procedure that given a set  $S \subseteq U$  as input, answers (correctly) if  $S \in \mathcal{I}$ .



Running time is the number of elementary operations + number of calls (i.e., queries) made to independence oracle. Let  $m = |U|$ . Then the running time is  $O(m \log m) + O(m \cdot \text{oracle}) = O(m \log m)$  which is polynomial-time, i.e., efficient algorithm.

### 4.3 Applications of Matroid Optimization

#### Maximum-Weight Bipartite Matching

##### bipartite

A graph  $G = (V, E)$  is called bipartite if  $V$  can be partitioned as  $L \cup R$  such that every edge has one end in  $L$  and one end in  $R$ . We call  $(L, R)$  bipartition of  $V$ .

##### matching

Given a graph  $G = (V, E)$ , a set  $M \subseteq E$  is called matching if  $|M \cap \delta(v)| \leq 1$  for all  $v \in V$ . I.e., for every  $v \in V$ , there is at most one edge of  $M$  incident to it.

##### Max-weight bipartite matching problem

Given a bipartite graph  $G = (V, E)$ , edge weights  $\{w_e\}_{e \in E}$ , find a matching  $M$  of maximum total edge weight.

Let  $L \cup R$  be bipartition of  $V$ . Define the following two matroid, having ground set  $E$ :

$$M_L = (U = E, \mathcal{I}_L = \{A \subseteq U : |A \cap \delta(v)| \leq 1, \forall v \in L\})$$

$$M_R = (U = E, \mathcal{I}_R = \{A \subseteq U : |A \cap \delta(v)| \leq 1, \forall v \in R\})$$

Because the graph is bipartite,  $\{\delta(v)\}_{v \in L}$  partitions  $E$ ,  $\{\delta(v)\}_{v \in R}$  partitions  $E$ . Then  $M_L$  and  $M_R$  are partition matroids. Thus  $A \subseteq E$  is a matching if and only if  $A \in \mathcal{I}_L \cap \mathcal{I}_R$ . Hence max-weight matching is equivalent to max-weight independent set in  $M_L$  and  $M_R$ . I.e., max-weight bipartite matching can be solved using matroid intersection.

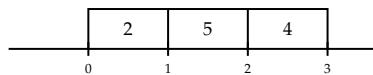
##### Scheduling Problem

Given a set  $U$  of jobs; each  $j \in U$  has **unit processing time**, deadline  $d_j \geq 1$  which is integer, weights  $w_j$  (can be arbitrary). We can only process one job at any time. Find a max-weight set  $S$  of jobs, and an ordering of  $S$  such that all jobs in  $S$  complete by their deadlines.

For example,

$j$	1	2	3	4	5	6
$d_j$	3	2	1	3	2	4
$w_j$	2	1.8	5.6	-0.4	-1	3

Let  $S = \{2, 5, 4\}$ : can order  $S$  so that all jobs complete by their deadlines.



Consider  $S = \{3, 2, 1, 4\}$ . We cannot order  $S$  so that all jobs complete by their deadlines since there are four jobs in  $S$  with deadlines  $\leq 3$ .

**schedulable**

Say that  $S \subseteq U$  is “schedulable” if there exists an order of  $S$  that completes all jobs by their deadlines.

**Exercise:**

If  $S$  is schedulable, then ordering jobs in  $S$  in increasing order of deadlines yields an ordering where all jobs complete by their deadlines.

$S$  is schedulable if and only if  $\forall t = 0, 1, \dots, |S|$ , (number of jobs in  $S$  with deadline  $\leq t$ )  $\leq t$

**Theorem 4.4**

$M = (U, \mathcal{I} = \{S \subseteq U : S \text{ is schedulable}\})$  is a matroid. Hence scheduling problem can be solved by solving MWIS for matroid  $M$ .

**Proof:**

Properties (a), (b) hold trivially. Consider exchange property. Let  $A, B \in \mathcal{I}$  with  $|A| < |B|$ . We have a notation: for  $S \subseteq U$ , let  $S_{\leq t} = \{j \in S : d_j \leq t\}$ . Consider smallest  $t \geq 0$  such that

$$|B_{\leq t'}| > |A_{\leq t'}| \quad \forall t' \geq t \quad (*)$$

Such a  $t$  exists, since for  $t = \max_{j \in B} d_j$ ,  $(*)$  holds and for  $t = 0$ ,  $(*)$  does not hold, and  $t \geq 1$ .

**Claim** There exists  $j \in B - A$  such that  $d_j = t$ .

If not, then all jobs in  $B$  with deadline equal to  $t$  are also in  $A$ . Then

$$\begin{aligned} |B_{\leq t-1}| &= |B_{\leq t}| - (\# \text{ jobs in } B \text{ with deadline} = t) \\ &> |A_{\leq t}| - (\# \text{ jobs in } B \text{ with deadline} = t) && \text{due to } (*) \\ &\geq |A_{\leq t}| - (\# \text{ jobs in } A \text{ with deadline} = t) \\ &= |A_{\leq t-1}| \end{aligned}$$

Contradicts  $t$  being the smallest value such that  $(*)$  holds.

**Claim**  $A' = A \cup \{j\}$  is schedulable where  $d_j = t$ .

Consider any  $t' = 0, 1, \dots, |A'|$ . Then

$$|A'_{\leq t'}| = \begin{cases} |A_{\leq t'}| \leq t' & \text{if } t' < t \\ |A_{\leq t'}| + 1 \leq |B_{\leq t'}| \leq t' & \text{if } t' \geq t \end{cases}$$

So by exercise,  $A'$  is schedulable. So  $M$  is a matroid. □

## Steiner Tree & Computational Complexity

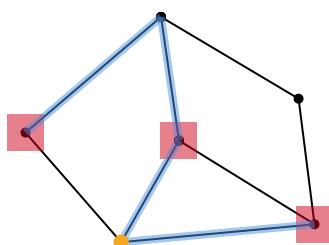
### 5.1 Minimum Steiner Tree Problem

#### Minimum Steiner Tree Problem

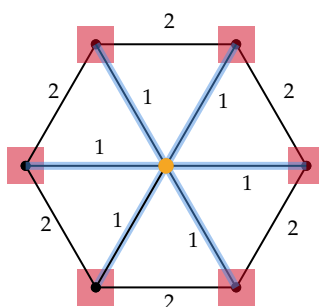
Given an undirected graph  $G = (V, E)$ , edge costs  $c_e \geq 0$ , and a set  $T \subseteq V$  called **terminals**, find a min-cost tree that spans (i.e., connects) all the vertices in  $T$ .

This is called Steiner tree for  $T$  and we will omit “for  $T$ ” if  $T$  is clear from the context.

Consider two examples.



Red boxes denote terminals  $T$ . Blue edges form a Steiner tree for  $T$ . Orange vertex is in Steiner tree but not in  $T$ . This is called a Steiner node. And we notice that this node is not necessary to include, but it might be useful in terms of the cost. For example,



Best Steiner tree containing only nodes from  $T$  has cost  $2(|T| - 1)$ . But blue edges form a Steiner tree of cost  $|T|$ .

**Useful Transformation** Let  $G = (V, E)$ .  $(G, c, T) \rightarrow (G', c', T)$  where  $G' = (V, E')$  is the complete graph on  $V$  and  $c'_{uv}$  is the shortest-path distance in  $G$  between  $u$  and  $v$ .

Note that  $\forall u, v, w \in V$ ,  $c'_{uw} \leq c'_{uv} + c'_{vw}$  ( $\triangle$ -ineq) since one  $u$ - $w$  path in  $G$  is shortest  $u$ - $v$  path in  $G$  + shortest  $v$ - $w$  path in  $G$ .

### Claim 5.1

- (a) Any Steiner tree  $F$  in  $G$  is also a Steiner tree in  $G'$ , and  $c'(F) \leq c(F)$ .
- (b) Any Steiner tree  $F'$  in  $G'$  yields a Steiner tree  $F''$  in  $G$  such that  $c(F'') \leq c'(F')$ .

### Proof:

For part (a), clearly  $F$  is a Steiner tree in  $G'$ . Also, note that  $\forall uv \in E$ ,  $c'_{uv} \leq c_{uv}$ , and so  $c'(F) \leq c(F)$ .

For part (b), take  $F'$ , and “expand” each edge  $uv \in F'$  to the shortest  $u$ - $v$  path  $P_{uv}$  in  $G$ , to get an edge-set  $\hat{F}$ . (So  $\hat{F} = \bigcup_{uv \in F'} P_{uv}$ , where  $P_{uv}$  is the shortest  $u$ - $v$  path in  $G$ .)

Clearly  $\hat{F}$  connects all vertices in  $G$ ; but it could contain cycles. We get  $F''$  from  $\hat{F}$  by simply removing edges from  $\hat{F}$  so that we get an acyclic graph. So

$$c(F'') \leq c(\hat{F}) \leq \sum_{uv \in F'} c(P_{uv}) = \sum_{uv \in F'} c'_{uv} = c'(F')$$

$\uparrow$   
 Since we removed  
 edges from  $\hat{F}$  to set  $F''$

□

$(G', c')$  defined as above is called **metric completion** of  $(G, c)$ .

By Claim 5.1, we can always solve Steiner tree problem on the metric completion of  $(G, c)$ . The optimum values do not change, and optimum solution in the metric completion gives us back an optimal solution in the original instance and vice versa.

## 5.2 MST-based algorithm

**Idea** Find an MST in  $G'[T] := (T, E'[T])$ <sup>1</sup> with respect to  $c'$ -costs and map this to a Steiner tree in  $G$  using Claim 5.1 (b).

---

### Algorithm 7: Algorithm for Minimum Steiner Tree Problems

---

- 1 Given instance  $(G, c, T)$ , consider the instance  $(G', c', T)$ , where  $(G', c')$  is the metric completion of  $(G, c)$ .
  - 2 Find an MST in  $G'[T] := (T, E'[T])$  with respect to  $c'$ -costs and map this back to a Steiner tree in  $G$  using Claim 5.1 (b).
- 

This is indeed a polytime algorithm, because all operations, including mapping back and so on, are all polynomial.

Cost of our solution  $\leq \text{MST}(G', c', T)$ , which  $c'$ -cost of MST in  $G'[T]$ . Let

$$\text{OPT} := \text{OPT}(G', c', T) = \text{OPT}(G, c, T)$$

where  $\text{OPT}(G', c', T)$  is  $c'$ -cost of an optimal solution for  $(G', c', T)$ .

### Theorem 5.2

(Cost of solution returned  $\leq$ )  $\text{MST}(G', c', T) \leq 2 \cdot \text{OPT}$ .

<sup>1</sup>subgraph of  $G'$  induced by  $T$ . It is also complete because  $G$  is complete

**Proof:**

Let  $F^*$  be optimal Steiner tree for  $(G', c', T)$ . Pick some  $r \in T$ , root  $F^*$  at  $r$ , and do a DFS traversal of  $F^*$  starting at  $r$ .

DFS traversal yields a tour  $Z$  (i.e., a closed walk or cycle with repeated nodes) that visits all vertices of  $F^*$ , and has  $c'$ -cost  $\leq 2c'(F^*) = 2 \cdot \text{OPT}$ .

Will “shortcut”  $Z$  to get a cycle  $\hat{Z}$  that visits every terminal exactly once. Suppose

$$Z : u_0 = r, u_1, u_2, \dots, u_{k-1}, u_k = u_0 = r$$

$$\hat{Z} : u_0 = r, u_{i_1}, u_{i_2}, \dots, u_{i_{\ell-1}}, u_{i_\ell} = u_0 = r$$

where for every  $j$ ,  $u_{i_j}$  is the first terminal after  $u_{i_{j-1}}$  in the tour  $Z$  that has not been visited before.

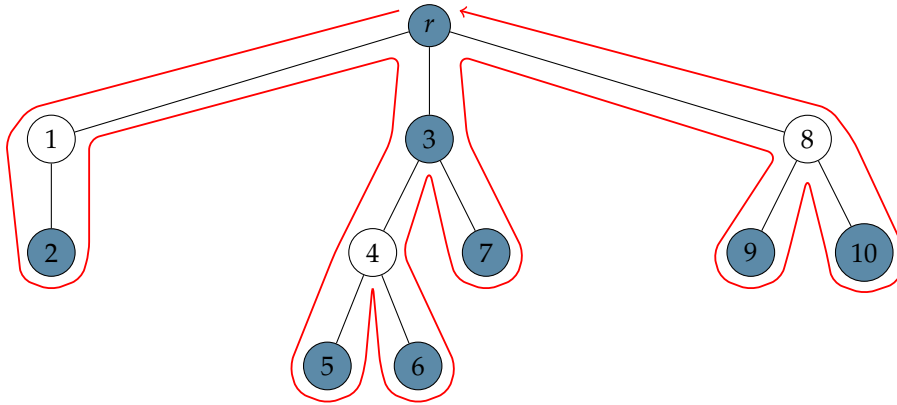
By  $\triangle$ -ineq, or the definition of  $c'$ , we have  $c'_{u_{i_j}u_{i_{j+1}}} \leq c'(Z_{u_{i_j}u_{i_{j+1}}})$  where  $Z_{u_{i_j}u_{i_{j+1}}}$  is the portion of  $Z$  between first occurrences of  $u_{i_j}$  and  $u_{i_{j+1}}$ .

Since  $Z_{u_{i_j}u_{i_{j+1}}}$  are disjoint for different  $j$  we have  $c'(\hat{Z}) \leq c'(Z)$ . Removing an edge from  $\hat{Z}$  gives a spanning tree in  $G'[T]$ , of cost  $\leq c'(\hat{Z}) \leq c'(Z) = 2 \cdot \text{OPT}$ .

Thus  $\text{MST}(G', c', T) \leq 2 \cdot \text{OPT}$ . □

**Example:**

Here is an example of how we get  $\hat{Z}$ .

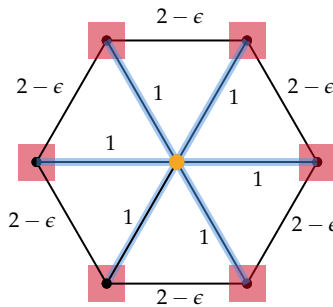


Terminal nodes are denoted by blue. So DFS traversal gives us  $Z$

$$r, 1, 2, 1, r, 3, 4, 5, 4, 6, 4, 3, 7, 3, r, 8, 9, 8, 10, 8, r$$

We mark first occurrences of terminals in  $Z$ . And these marked vertices form  $\hat{Z}$ .

Were we dumb in analyzing this algorithm? Actually no. Consider an example.



where  $\epsilon > 0$  is small. Our algorithm returns all edges of cost  $2\epsilon$ , thus the total cost is  $(2 - \epsilon)(|T| - 1)$ .

But the optimal solution is blue edges, of cost  $|T|$ . We see the factor is approximately 2. Can we design a better algorithm? For this problem, there is no known efficient algorithm that is always guaranteed to return an optimal solution. Moreover, researchers don't believe such an algorithm exists.

#### $\alpha$ -approximation algorithm

An  $\alpha$ -approximation algorithm, where  $\alpha \geq 1$  is the approximation factor, for a minimization problem, is an efficient (i.e., polytime) algorithm that on every instance returns a solution of cost  $\leq \alpha \cdot \text{OPT}$ .

So above algorithm for Steiner tree is a 2-approximation algorithm for the Steiner tree problem.

## 5.3 Computational Complexity

$P$  is class/set of all problems that can be solved by polytime algorithms. For example, Shortest Paths, MST, MCA  $\in P$ . Other YES/NO problems:

- **IsComposite**: Given integer  $n \geq 2$ , is  $n$  a composite number; i.e, do exist integers  $x, y \geq 2$  such that  $n = xy$ ?
- **Factoring**: Given integer  $n \geq 2$ , decide if  $n$  is composite and if so, find integers  $x, y \geq 2$  such that  $n = xy$ .
- **IsPrime**: Given integer  $n \geq 2$ , is  $n$  a prime number?

---

#### Algorithm 8: Factoring algorithm: Sieve of Eratosthenes

---

- 1 Consider all integers  $x$  such that  $2 \leq x \leq \sqrt{n}$ , and see if  $x$  divides  $n$ .
  - 2 If so, then  $n$  is composite and  $x, \frac{n}{x}$  is a factorization of  $n$ .
- 

Is the above a polytime algorithm? No, since number of iterations  $\approx \sqrt{n}$ , which is not polynomial in  $O(\log n)$  which is the input size.

Open question: Is **Factoring**  $\in P$ ? We already know that **IsPrime**  $\in P$ , which is only settled in around 2003. This implies **IsComposite**  $\in P$ .

Decision problem is problem with a YES/NO answer. We can take any optimization problem and consider a decision-version of the problem:

- Is there a solution of cost  $\leq k$  (for minimization problems)
- Is there a solution of value  $\geq k$  (for maximization problems)

Here is decision version of MST (**DecMST**): given  $(G, \{c_e\})$ , and a number  $k$ , is there a spanning tree of cost  $\leq k$ ? If **MST**  $\in P$ , then **DecMST**  $\in P$ .

## 5.4 Polynomial-Time Reductions

#### polytime reduction

Given problems  $A, B$ , we say  $A$  **reduces in polynomial time to  $B$** , or  $A$  is **polytime reducible to  $B$** , denote  $A \leq_p B$ , if we can solve problem  $A$  using a polynomial number of elementary operations + a polynomial number of calls to an algorithm for problem  $B$ .

$MST \leq_p MCA$ . Given input  $(G = (V, E))$  undirected graph,  $\{c_e\}_{e \in E}$  to MST problem, create  $(\overset{\leftrightarrow}{G}, \overset{\leftrightarrow}{c})$  the bidirected version of  $G$ , by creating edges  $(u, v), (v, u)$  in  $\overset{\leftrightarrow}{G}$  for every edge  $uv$  of  $G$ , and give these two edges cost  $c_{uv}$ . Solve MCA on  $(\overset{\leftrightarrow}{G}, \overset{\leftrightarrow}{c}, r)$  where  $r$  is any node of  $V$ , and map back MCA to undirected edge of  $G$  to get an MST.

$MST \leq_p$  Min Steiner Tree Problem. This is true because MST is a special case of Min Steiner Problem, when edge costs are  $\geq 0$ . For MST, can always ensure (i.e., move to) nonnegative edge costs by adding a suitable large constant to all edge costs.

Suppose  $A \leq_p B$ . Then

1. If  $B \in P$ , then  $A \in P$ .
2. Equivalently, if  $A \notin P$ , then  $B \notin P$ .

Consider

- SORTING: Given  $n$  numbers  $a_1, \dots, a_n$ , sort them in an increasing order.
- MIN: Given  $n$  numbers  $a_1, \dots, a_n$ , find the minimum of these  $n$  numbers.

$\text{SORTING} \leq_p \text{MIN}$ : Can sort using  $n$  calls to an algorithm for MIN. Also  $\text{MIN} \leq_p \text{SORTING}$ .

Consider <sup>2</sup>

- LPFeas: Given  $A \in \mathbb{Q}^{m \times n}$ , is the system  $Ax \leq b, x \geq 0$  feasible?
- LP0pt: Given  $A \in \mathbb{Q}^{m \times n}, b \in \mathbb{Q}^m, c \in \mathbb{Q}^n$ , does the LP

$$\begin{aligned} \max \quad & c^T x \\ \text{s.t.} \quad & Ax \leq b \\ & x \geq 0 \end{aligned}$$

have an optimal solution?

$\text{LPFeas} \leq_p \text{LP0pt}$ . We can set  $c = 0$ , and call algorithm for LP0pt. Using duality we can show  $\text{LP0pt} \leq_p \text{LPFeas}$ .

Decision problem LPImprov: Given  $A \in \mathbb{Q}^{m \times n}, b \in \mathbb{Q}^m, c \in \mathbb{Q}^n$ , and  $z \in \mathbb{Q}$ , determine if there is a feasible solution to the LP:

$$\begin{aligned} \max \quad & c^T x \\ \text{s.t.} \quad & Ax \leq b \\ & x \geq 0 \end{aligned}$$

with objective value  $> z$ ? Using duality, can show that  $\text{LPImprov} \leq_p \text{LPFeas}$ .

Note that simplex method for solving LPs is NOT a polytime algorithm.

Consider

- MST: Given connected graph  $G = (V, E)$ , edge costs  $\{c_e\}$  that are positive integers, find cost of MST.
- Decision version of MST, DecMST: Given connected  $G = (V, E)$ , positive integer edge costs  $\{c_e\}$ , and an integer  $k > 0$ , is there a spanning tree of cost  $\leq k$ ?

We observe that  $\text{DecMST} \leq_p \text{MST}, \text{MST} \leq_p \text{DecMST}$ .

Given MST instance  $(G, c)$ , we can deduce that every spanning tree has cost  $\geq LB := 1$ , and MST has cost  $\leq UB := (n - 1)c_{\max}$ , where  $c_{\max} := \max_e c_e$ .

<sup>2</sup>Here we choose to work with rational numbers because it's hard to specify the input size for irrational numbers

**Algorithm** Start at  $k = UB$ . Solve DecMST with this  $k$ ; if answer is YES, set  $k \leftarrow k - 1$  and REPEAT.

This is not a polytime reduction, since we could be calling algorithm for DecMST  $\approx UB$  times, but  $UB \approx (n - 1)c_{\max}$  is not polynomially bounded in input size.

We can get polynomial reduction by doing binary search. Binary search calls algorithm DecMST  $O(\log_2 UB)$  which is polynomially bounded in input size.

It's not hard to see that  $\leq_p$  is transitive.

Consider two problems:

- MST: find cost of MST
- A-MST: Find a min-cost spanning tree.

We can see that  $MST \leq_p A\text{-}MST$ .  $A\text{-}MST \leq_p MST$ : Let  $k = MST(G, c)$ . Let  $e$  be some edge of  $G$ . Then  $k' = MST(G - e, \{c_f\}_{f \in G - e})$ . If  $k' > k$ , then  $e$  belongs to every MST. If  $k' \leq k$ , then  $G \leftarrow G - e$ ,  $c \leftarrow \{c_f\}_{f \neq e}$ .



## P & NP

---

Consider decision version of Steiner tree:

**DecSteiner:** Given an instance  $(G, c, T)$  and a number  $k$ , is there a Steiner tree of cost  $\leq k$ ?

**Observation** If DecSteiner instance is a YES instance, then there is a simple certificate (A steiner tree of cost  $\leq k$ ) to convince one that answer is indeed YES.

### verifier/certifier

A verifier/certifier  $V$  for a decision problem  $\Pi$  is an algorithm that takes two inputs,  $x$ : an instance of  $\Pi$ , and a “certificate”  $y$ , and outputs YES or NO. It satisfies

- If  $x$  is a YES instance, then there exists  $y$  such that  $V(x, y) = \text{YES}$ .
- If  $x$  is a NO instance, then for all  $y$ ,  $V(x, y) = \text{NO}$ .

### 6.1 NP

**NP** a decision problem  $\Pi$  is in the class **NP** if there exists a polynomial  $p$  and a polytime verifier  $V$  for  $\Pi$  such that

- For every YES instance  $x$  of  $\Pi$ , there exists a certificate  $y$  with  $\text{size}(y) \leq p(\text{size}(x))$  such that  $V(x, y) = \text{YES}$ .
- (For every NO instance  $x$ ,  $V(x, y) = \text{NO}$  for all  $y$ )

Informally, a decision problem is in **NP** if its YES instances have “short”, efficiently verifiable certificates.

#### Example:

**DecSteiner**  $\in$  **NP**. It has certificate  $y$ : Steiner tree of cost  $\leq k$ . Verifier  $V$ : check  $y$  is indeed a Steiner tree of cost  $\leq k$ .

#### Example:

Recall **IsComposite**. It is in **NP**. Certificate  $y$ : a factor of  $n$  where  $2 \leq y < n$ . Verifier: check  $y$  divides  $n$  (and  $2 \leq y < n$ ).

**Example:**

IsPrime  $\in NP$ ? Yes, certificate for YES instance uses results in number theory, Fermat's little theorem. We know that IsPrime  $\in P$ .

**Claim 6.1**

If  $\Pi$  is a decision problem in  $P$ , then  $\Pi$  is also  $NP$ . In other words,  $P \subseteq NP$ .

**Proof:**

Let  $A$ : polytime algorithm for  $\Pi$ . Then, we can construct a polytime verifier  $V$  for  $\Pi$  as follows: on input  $(x, y)$ ,  $V$  returns  $A(x)$ .  $\square$

**Exercise:**

Prove that  $P \neq NP$ .

**Example:**

3-SAT: Given a boolean formula  $F$  of the form

$$F = C_1 \wedge C_2 \wedge \cdots \wedge C_m,$$

where each  $C_i$  is a clause of the form  $y_{i_1} \vee y_{i_2} \vee y_{i_3}$  where each  $y_i$  (called literal) is  $x_i$  or  $\bar{x}_i$ .

Is there a True/False assignment of the variable under which  $F$  evaluates to T? Such an assignment is called a satisfying assignment.

Not hard to see 3-SAT  $\in NP$ : certificate  $y$ : satisfying assignment. Verifier: check that  $y$  is indeed a satisfying.

**Example:**

Set-Cover: Input: a universe  $U$ , and a collection  $\mathcal{S} = \{S_1, \dots, S_n\}$  of subsets of  $U$ , and an integer  $k \geq 0$ .

Are there  $k$  subsets  $S_{i_1}, \dots, S_{i_k} \in \mathcal{S}$  such that  $(S_{i_1} \cup \cdots \cup S_{i_k}) = U$ ?

Set-Cover  $\in NP$ : a certificate for YES instance is simply  $k$  sets of  $\mathcal{S}$  whose union is  $U$ .

## 6.2 NP-hard and NP-complete

### NP-hard, NP-complete

A problem  $B$  is called:

- NP-hard: if  $X \leq_p B \forall x \in NP$
- NP-complete: if  $B \in NP$  and  $B$  is NP-hard (i.e.,  $x \leq_p B \forall x \in NP$ )

Intuitively, NPC are hardest problems in class  $NP$ .

Suppose  $Y$  is a NPC problem. Then

- if  $Y \in P$ , then  $P = NP$ . Since  $X \leq_p pY$  for all  $X \in NP$  and  $Y \in P$ , we get that  $x \in P \forall x \in NP$ .
- If  $Y \notin P$ , then  $P \neq NP$ .

How do we show that a problem  $Y$  is NP-hard or NPC?

- Pick a suitable NP-hard problem  $B$ .

- Show  $B \leq_p Y$ . Then  $Y$  is NP-hard by transitivity of  $\leq_p$ .
- If we want to show that  $Y$  is NP-complete, then also SHOW  $Y \in NP$ , which is usually easy.

**Note:**

If  $Y'$  is the decision version of optimization problem  $Y$  and  $Y'$  is NP-hard, then  $Y$  is also NP-hard.  
Since  $Y' \leq_p Y$ .