

Lecture 6

September 25, 2023

Instructor: Sepehr Assadi

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

Topics of this Lecture

1	Minimum Spanning Trees	1
1.1	The Basics Rules of MSTs	2
2	Refresher on the Classical Algorithms for MSTs	2
2.1	Kruskal's Algorithm	3
2.2	Prim's Algorithm	4
2.3	Boruvka's Algorithm	5
3	Fredman and Tarjan's Algorithm in $O(m \log^*(n))$ Time	6
3.1	Warm-Up: Improving Prim's Runtime via Fibonacci Heaps	6
3.2	The Fredman-Tarjan Algorithm	6

1 Minimum Spanning Trees

We now begin the second half of the course, starting with the problem of finding a **minimum spanning tree (MST)** of *weighted* graphs. You have most likely seen this course in your previous algorithms course and you can also refresh your background on this problem by checking Erikson's *Algorithms* book (or any other textbook on algorithms that you like). We will provide a quick review in the following.

Let $G = (V, E)$ be an undirected *connected* graph with *positive* weights $w : E \rightarrow \mathbb{R}^+$ on its edges. Recall that a spanning tree of a connected subgraph G is any subgraph of G which is a tree – a tree itself is a connected subgraph which has no cycle. With a slight abuse of notation, for any tree T , we write $w(T) := \sum_{e \in T} w(e)$ to denote the sum of the weights of edges in T .

Problem 1 (Minimum Spanning Tree (MST)). The minimum spanning tree (MST) problem is defined as follows: Given a graph $G = (V, E)$ and positive edge-weights $w : E \rightarrow \mathbb{R}^+$, find a spanning tree T of G with minimum weight $w(T)$.

Distinct-weights assumption: We can assume without loss of generality that the edge weights $w(e)$ for $e \in E$ are *distinct*. This can be achieved for instance by using a consistent tie-breaking rule using the IDs of the edges. A standard consequence of this assumption is that the MST of any given graph G with distinct weights is *unique*.

Proposition 1. *In any graph $G = (V, E)$ with weights $w : E \rightarrow \mathbb{R}^+$, if the weights are distinct, then the MST of G is unique.*

Proof. Proof by contradiction: suppose there are two different MSTs T_1 and T_2 for G . Let e_1 be any edge in T_1 not in T_2 . Suppose we add e_1 to T_2 . This will create a cycle C in $T_2 \cup \{e_1\}$. We consider two cases:

- *Case 1:* There is at least one edge $e_2 \neq e_1$ in C with $w(e_2) > w(e_1)$: Remove e_2 from T_2 and instead add e_1 . Let $T := T_2 \cup \{e_1\} \setminus \{e_2\}$. This is also a spanning tree of G but now $w(T) < w(T_2)$ as we removed a “heavy” edge and brought in a “lighter” one. This is a contradiction with T_2 being a MST so this case cannot happen.
- *Case 2:* Every edge in the cycle C has weight higher than e_1 . Remove e_1 from T_1 to create two connected components S and $V \setminus S$. Because C was a cycle with e_1 , there is still a path in C from the endpoints of e_1 without using this edge and that path has at least one edge, say, f , between S and $V \setminus S$. Let $T := T_1 \setminus \{e_1\} \cup \{f\}$. This is also a spanning tree of G but now $w(T) < w(T_1)$, a contradiction with T_1 being a MST, so this cannot happen either.

Given the edge-weights are distinct, the only possible cases are the above and we showed neither can happen, hence proving the statement by contradiction. \square

Consequently, throughout the course, unless explicitly stated otherwise, we assume the edge-weights are distinct and thus, by **Proposition 1**, **MST of the input graph is unique**. In particular, we can refer to it as *the* MST of G (instead of a/some MST of G) which makes the analysis (mostly its exposition) easier.

1.1 The Basics Rules of MSTs

There are two general *rules* that are used quite frequently in the design of essentially every MST algorithm. You have most likely seen these rules before but perhaps not in this exact formulation. You may want to refer to **Figure 1** for a simple illustration of the proofs of these rules.

Cut Rule: *In any graph $G = (V, E)$, the minimum weight edge e in any cut S always belongs to the MST of G . This rule allows us to determine which edges to include in the MST.*

Proof. Suppose, by contradiction, that $S \subseteq V$ is a cut in G such that its minimum weight edge e is not part of the MST T . Add e to T which creates a cycle C and this cycle should have another edge f crossing the cut S . Consider $T' := T \cup \{e\} \setminus \{f\}$, which is still a spanning tree of G but has $w(T') < w(T)$ as $w(e) < w(f)$ since e is the minimum weight edge of the cut S . This is a contradiction with T being the MST of G . \square

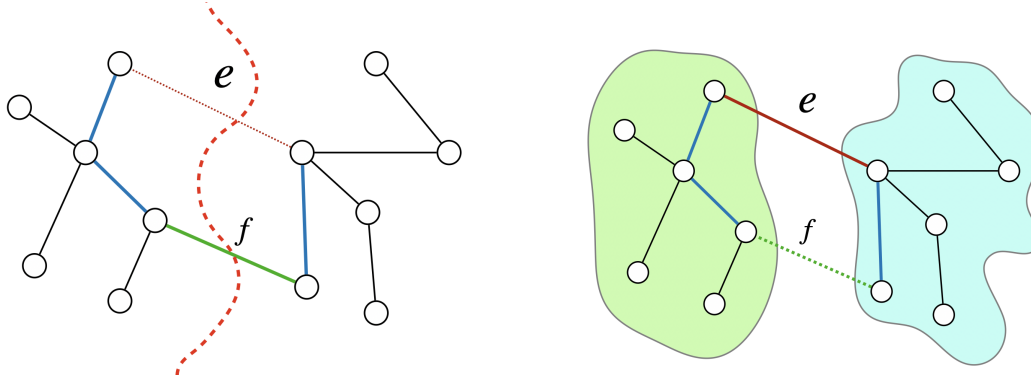
Cycle Rule: *In any graph $G = (V, E)$, the maximum weight edge e in any cycle C never belongs to the MST of G . This rule allows us to determine which edges to exclude from the MST.*

Proof. Suppose, by contradiction, that C is a cycle in G such that its maximum weight edge e is part of the MST T . Remove e from T to get two connected components, and find an edge $f \neq e$ in C that connects back these components. Consider $T' := T \setminus \{e\} \cup \{f\}$, which is still a spanning tree of G . We have $w(T') < w(T)$ as $w(f) < w(e)$ since e is the maximum weight edge of cycle C . This contradicts T being the MST. \square

2 Refresher on the Classical Algorithms for MSTs

This section reviews classical algorithms for MSTs—*Kruskal’s*, *Prim’s*, and *Boruvka’s* algorithms—that you most likely have seen in your previous courses (at least some subset of them¹). You can safely skip this section and continue to the next one if you are already familiar with these algorithms.

¹Sadly, Boruvka’s algorithm seems to be missing from many standard textbooks on algorithms even though I think it is actually simpler than both the other two alternatives!



(a) Here, e is the minimum weight edge of the cut shown by red dashed line and thus should replace f in the MST. The cycle C is the thick (blue, green) edges and the dashed line for e .

(b) Here, e is the maximum weight edge of the cycle shown by the thick (blue, green, red) edges. Removing e from the MST creates two connected components that should be connected by the edge f instead.

Figure 1: An illustration of the proofs of **cut rule** and **cycle rule**.

2.1 Kruskal's Algorithm

To state Kruskal's algorithm, we first need to review the *union-find* data structure.

Definition 2. A **union-find** is a data structure supporting these operations on input set $\{1, 2, \dots, n\}$:

- **initialize**: starts the data structure by creating n *disjoint* sets S_1, \dots, S_n where $S_i := \{i\}$;
- **union**(i, j): Given index of two sets i, j , replace S_i and S_j with $S_i \cup S_j$ and the index of the new set with $\min \{i, j\}$;
- **find**(e): given an element $e \in [n]$, returns the index of the set S_j where $e \in S_j$.

There are different ways of implementing a union-find data structure and this problem was studied extensively in 60's, 70's, and 80's and at this point is extremely well-understood. We will not go over these implementations in this course². For our purpose, we can use the simplest (non-trivial) implementation of this data structure that implements **initialize** in $O(n)$ time and **union** and **find** in $O(\log n)$ time.

Kruskal's algorithm works by iterating over the edges in the order of their weights and include every edge that does not create a cycle; this is done by maintaining the connected components in a union-find data structure and only pick an edge if it connects two different components and merge those components.

Algorithm 1 (Kruskal's Algorithm).

- (i) Create a union-find data structure D on n vertices in V by running **initialize**. Let $T = \emptyset$. Sort the edges of G in increasing order of their weights.
- (ii) For $i = 1$ to m in the sorted ordering of edges:
 - (a) Let $e_i = (u_i, v_i)$. Let $a = D.\text{find}(u_i)$ and $b = D.\text{find}(v_i)$.
 - (b) If $a = b$ continue to the next edge; otherwise, $T \leftarrow T \cup \{e_i\}$ and run $D.\text{union}(a, b)$.
- (iii) Return T .

²But you are encouraged to check this data structure online for mind boggling aspects of these implementations

Proof of Correctness: The fact that [Algorithm 1](#) returns a spanning tree is immediate: any edge is ignored only if it creates a cycle (as it is inside a single connected component) and thus non-ignored edges should maintain the connectivity of the input graph. To see that T is the MST of G we can apply the **cycle rule**: any edge e ignored by the algorithm has both its endpoints inside a single connected component and thus forms a cycle inside that component; given the sortedness of the edges, this ignored edge is the maximum weight of this cycle, and thus by the **cycle rule** cannot be part of the MST. So, any ignored edge cannot be part of the MST and thus the non-ignored edges T should form the MST.

Runtime Analysis: The first line of the algorithm takes $O(m \log m) = O(m \log n)$ time (as $n \leq m \leq n^2$ and so $\log n = \Theta(\log m)$) to sort all the edges (using, say, merge sort). Initializing the union-find takes $O(n)$ time and each of its operations takes $O(\log n)$ time, so each iteration of the for-loop also takes $O(\log n)$ time, making the total runtime of the algorithm $O(m \log n)$. Finally, it is worth pointing out that we can implement for-loop *much faster* almost (but not quietly) in linear time; however that will not reduce the runtime as we still need to sort the edges in $O(m \log m)$ time in general.

2.2 Prim's Algorithm

We now consider another textbook algorithm for MSTs, the *Prim's Algorithm*. To do so, we need to recall the definition of *priority queue* data structures.

Definition 3. A **priority queue** is a data structure Q supporting the following operations on a universe \mathcal{U} of elements:

- **insert**(e, v): Given an element $e \in \mathcal{U}$ and a value $v \in \mathbb{R}$, inserts the element e with value v in Q ;
- **decrease-key**(e, v): Given an element $e \in \mathcal{U}$ and a value $v \in \mathbb{R}$, changes the value of e to *minimum* of v and its old value;
- **extract-min**(): Removes the element e with the smallest value v from H and return (e, v) .

The Kruskal's algorithm is about maintaining *multiple* connected component of the forest F simultaneously and updating them (merge them together in the algorithm). On the other hand, Prim's algorithm only contains *one non-trivial* connected component (all other components are singleton vertices) which we "grow" until it contains the entire graph. The strategy for growing the component is also very basic: we simply pick the edge with minimum weight that goes out of this component in every step.

Algorithm 2 (Prim's Algorithm).

- (i) Pick an *arbitrary* vertex s of the graph. Let T be a subgraph of G consisting of only the vertex s for now and Q be a priority queue with keys being vertices in V and key-values being an edge $edge(v)$ and its weight (used for comparisons in **decrease-key** and **extract-min** operations).
- (ii) For every vertex $v \in N(s)$ set $edge(v) = (s, v)$ and run $Q.\text{insert}(v, w(edge(v)))$.
- (iii) For $n - 1$ steps:
 - (a) Let v be the vertex with the smallest value in Q by running $Q.\text{extract-min}$ and $edge(v)$ be the edge of v with weight equal to value of v in Q . Add vertex v and the edge $edge(v)$ to T .
 - (b) For any vertex $u \in N(v)$, if $u \notin T$ update $Q.\text{decrease-key}(u, w(uv))$ (if u is not in Q , we instead run $Q.\text{insert}(u, w(uv))$).
- (iv) Return T as the MST of G .

Proof of Correctness: There are two steps in proving the correctness of Prim’s algorithm: (1) it outputs a spanning tree, and (2) the weight of this spanning tree is minimum, namely, it is an MST. The proof of first part is simply as follows: the graph T created by the algorithm consists of $n - 1$ edges (given the $n - 1$ steps of the for-loop), and does not have a cycle as every edge added is incident on a vertex of degree v at the time (the vertex included in T). Thus T is a spanning tree of the input.

To prove T is an MST we apply the **cut rule**. Let e be an edge added to T by the algorithm. Let C be the connected component of T before adding the edge e . Since all vertices in C have been marked, all edges incident on them have been added to Q . This means that all edges between T and $N(T)$, the set of neighbors of T , have been included in Q (by way of running $Q.\text{decrease-key}$). Since e has the minimum weight inside Q , this means that it has the minimum weight among all edges between T and $N(T)$. Thus, e is the minimum weight edge of the cut T and hence, by the **cut rule** e belongs to the MST. This proves all $n - 1$ edges inserted to T belong to the MST, hence T is the MST of the input graph.

Runtime Analysis: The algorithm requires $O(n)$ **insert** operation, $O(m)$ **decrease-key**, and $O(n)$ **extract-min** operations (and additionally takes $O(m + n)$ time for standard bookkeeping). Thus, if we implement the priority queue via a *min-heap* (a.k.a. a *binary heap*) that runs in each operation in $O(\log |Q|)$ time, we get that the total runtime of the algorithm is $O(m \log n)$.

2.3 Boruvka’s Algorithm

Finally, we get to Boruvka’s algorithm, which, as far as I know, is the oldest algorithm for MSTs and in my opinion is also simplest as it does not require even any particular data structure. It works by start having each vertex its own connected component and then proceeds in *rounds*. In each round, the algorithm picks the minimum weight edge going out of each component, add them to the tree, and merge those components to create larger connected components. Then, it goes to the next round and continues this until the whole graph is a single connected component.

Algorithm 3 (Boruvka’s Algorithm).

- (i) Start by creating a set $S(v) = \{v\}$ for each vertex $v \in V$ and an empty subgraph $T = \emptyset$ initially.
- (ii) In each **round**:
 - (a) For each set S , find the minimum weight edge going out of S , i.e., $\min \{w(u, v) \mid u \in S, v \notin S\}$. Let e_S be the edge computed for the set S .
 - (b) Add every edge $e_S = (u, v)$ to the subgraph T and merge the sets of $S(u)$ and $S(v)$ into a single set S .
- (iii) Repeat the rounds until all vertices are merged into a single set and return T .

Proof of Correctness: As before, we first argue that T is a tree. Given that the maintained sets in the algorithm are connected components of T , the subgraph T is connected. Any edge added by the algorithm in a round also connects two different connected components. The only concern is that what if the edges added in the same round form a cycle? But, this cannot happen because the weight of the picked edges in any path should be strictly decreasing (as in, if u picks an edge to v , but v picks an edge to w instead, then $w(u, v) > w(v, w)$) and thus the edges cannot form a cycle. As such, T is indeed a tree.

To argue T is the MST of G , we again apply the **cut rule**: any edge inserted by the algorithm is the minimum weight edge going out of its connected component (hence the cut) and thus by the **cut rule** belongs to the MST.

Runtime Analysis: Each round of the algorithm takes $O(m)$ time as it simply involves going over all edges of the graph twice (once from each endpoint). The number of rounds is also $O(\log n)$ as we argue next.

In each round, each connected components gets merged with at least another one, and thus the number of components decreases by, at least, a factor of two. Hence, after at most $\log n$ rounds, only one components remain and the algorithm terminates. Hence, the total runtime is $O(m \log n)$.

This concludes our summary of the background on classical algorithms for MST. You can find a lot more information in Erikson’s *Algorithms* book or any other textbook on undergraduate-level algorithms.

3 Fredman and Tarjan’s Algorithm in $O(m \log^*(n))$ Time

We now go over one of the earliest “advanced” algorithms for MST presented by Fredman and Tarjan in [FT87] as part of—and one key application of—their introduction of *Fibonacci heaps*. The Fredman-Tarjan algorithm runs in $O(m \log^*(n))$ time, where $\log^*(n)$ is the number of times we need to take $\log(\cdot)$ until the number reaches one³. The algorithm is quite simple and even more so elegant and beautiful.

Let us start by defining Fibonacci heaps that implement priority queues (Definition 3) faster than a typical min-heap.

Proposition 4 ([FT87]). *There is a data structure H , called a **Fibonacci heap**, that implements a priority queue (Definition 3) such that **insert** and **decrease-key** operations takes $O(1)$ and **extract-min** operation takes $O(\log |H|)$ amortized time where $|H|$ denotes the maximum number of elements in the heap during the entire execution.*

The construction of Fibonacci heaps is not particularly simple and we will not go over that in this class. Instead, we show how to use this data structure to get faster MST algorithms.

3.1 Warm-Up: Improving Prim’s Runtime via Fibonacci Heaps

Using Fibonacci heaps in Prim’s algorithm in place of its min-heap immediately leads to improving its runtime. As stated, Prim’s algorithm involves $O(n)$ **insert** operations, one for each vertex of the graph, $O(m)$ **decrease-key** operations, one each time we visit an edge of the graph, and $O(n)$ **extract-min** operation, one each time we mark a new vertex. Thus, using Fibonacci heaps, by Proposition 4, the runtime of Prim’s algorithm actually reduces to $O(m + n \log n)$ time.

As long as the input graph is even slightly dense, i.e., $m = \Omega(n \log n)$ edges, Prim’s algorithm equipped with Fibonacci heaps runs in linear-time in the input size which is optimal. But, when m is $\Theta(n)$ itself, then this algorithm does not provide any improvement over the runtime of the original Prim’s algorithm.

3.2 The Fredman-Tarjan Algorithm

We next go over Fredman-Tarjan algorithm that improves this runtime for all values of m and prove the following theorem.

Theorem 5 ([FT87]). *There is a deterministic algorithm for the minimum spanning tree problem that runs in $O(m \log^*(n))$ time.*

The key observation behind the Fredman-Tarjan algorithm is that the runtime of Prim’s algorithm with a Fibonacci heap depends on the size of the heap, which itself is the equal to the *frontier* of the graph search, namely, vertices waiting in the priority queue to be examined next. So, instead of growing the frontier to become “too much to handle”, we can limit its size and instead run Prim’s algorithm for another starting vertex and keep doing this – one then needs to recombine these partial trees recursively.

We now state the formal algorithm.

³This is an extremely slow growing function and even if we take n to be the number of atoms in the entire universe, $\log^*(n)$ will be at most 6 still!

Algorithm 4 (Fredman-Tarjan Algorithm [FT87]).

- The algorithm runs in **rounds**. In round i , it works with a graph $G_i = (V_i, E_i)$ with n_i vertices and m_i edges which is obtained from G by *contracting* some vertices chosen in previous rounds; we take $G_1 = G$ at the beginning. In round i , it also use a *threshold* t_i on the size of the Fibonacci heaps it uses which we set to be $t_i := 2^{2m/n_i}$ (note that the nominator in the exponent is m independent of the round but the denominator depends on the round).
- In each round i , the algorithm does as follows:

- (i) Start with every vertex $v \in V_i$ being unmarked.
- (ii) Pick an arbitrary unmarked vertex and run Prim's algorithm from it to get a tree T_v . In particular, use Fibonacci heaps on the neighborhood of T_v , defined as

$$N(T_v) := \{u \in V_i \setminus T_v \mid \exists w \in T_v \text{ such that } u \in N(w)\};$$

in each step, take the vertex with minimum edge-weight from $N(T_v)$ (via **extract-min** operation of the heap) and expand T_v from that vertex (note that vertices in $N(T_v)$ can be marked already but we still include them in $N(T_v)$).

- (iii) If at the end of each step of Prim's algorithm in the previous line, $|N(T_v)| \geq t_i$, or that the next vertex to explore is already a marked vertex (i.e., T_v just got an edge to an already marked vertex), then terminate this branch of Prim's algorithm, mark all vertices in T_v , and go to line (ii). Continue as long as there is any unmarked vertex left.
- (iv) Let C_1, C_2, C_3, \dots be the connected components created by different T_v 's picked in the above steps. Create G_{i+1} for the next round by contracting these components in G_i into a single vertex. Then go to the next round and continue this until the whole graph is contracted to a single vertex and we have a single spanning tree (obtained by combining the trees computed in different steps).

Proof of Correctness. The proof of correctness of the algorithm is pretty simple and follows from the correctness of Prim's and Boruvka's algorithms. We first have that **Algorithm 4** always finds a spanning tree of the input graph the same way Boruvka's algorithm always finds a spanning tree. We then have that each edge e added to the tree T is a minimum weight edge of some cut and thus we can apply the **cut rule** to argue this edge belongs to the MST. The cut itself is found by considering the connected component of G on the set of vertices in T (and possibly by un-contracting the vertices in rounds > 1), the same exact way as in Prim's algorithm.

Runtime analysis. We now switch to the runtime which is the main part of the analysis for this new algorithm. Each round involves visiting each edge at most twice (from either side) and inserting/updating their weights into a Fibonacci heaps; as well as removing each vertex from the heap at most once via the **extract-min** operation (after that, the vertex is marked one way or another and will not be inserted to the heap again). As size of each heap is capped at t_i by the construction of the algorithm, the total runtime of the round i , by **Proposition 4**, is $O(m_i + n_i \cdot \log t_i)$ time. Given the choice of $t_i := 2^{2m/n_i}$, this total runtime is $O(m_i + m)$ time. Finally, as $m_i \leq m$, we have that **Algorithm 4** takes $O(m)$ time in each round.

Consequently, to bound the total runtime of the algorithm, we only need to bound the number of rounds before the algorithm contract the entire graph into a single connected component. This is done using the following key lemma, that shows the thresholds are growing exponentially in each round.

Lemma 6. *For every round i of **Algorithm 4**, we have $t_{i+1} \geq 2^{t_i}$.*

Proof. We are going to bound n_{i+1} as a function of t_i and use the fact that $t_{i+1} = 2^{2m/n_{i+1}}$ to conclude the proof. By definition, n_{i+1} is equal to the number of connected components found at the end of the round.

While we are truncating Prim's algorithm to make its heap small, we still do not expect to have “very small” components using the threshold we are using.

In particular, we claim that, for every component C in round i :

$$\sum_{v \in C} \deg_{G_i}(v) \geq t_i; \quad (1)$$

namely, sum of the degrees of vertices (in G_i) in a component C is at least equal the threshold. The proof of Eq (1) is as follows. For each sub-tree T created in the algorithm in round i and each vertex v extracted from the heap (i.e., marked) during its construction, we put at most $\deg_{G_i}(v)$ vertices in the heap corresponding to T (we may insert less if those vertices are already inserted). This means that size of the heap can only increase to sum of the degrees of extracted/inserted vertices. Thus, if $|N(T)| \geq t_i$, we know that all vertices in T will be marked and so collectively they have contributed at least t_i to the sum of the degrees, hence proving Eq (1). On the other hand, if we terminate the search because we connected to another marked vertex, then, by the prior argument (or rather inductively) sum of the degrees in that component was already t_i and we are only adding these vertices to that component which can only increase the total degrees. This proves Eq (1).

Now, let $C_1, \dots, C_{n_{i+1}}$ be the connected components found in this step (recall that n_{i+1} is equal to the number of these connected components). We have,

$$\begin{aligned} 2m_i &= \sum_{v \in V_i} \deg_{G_i}(v) && \text{(by the handshaking lemma)} \\ &= \sum_{j=1}^{n_{i+1}} \sum_{v \in C_j} \deg_{G_i}(v) \\ &\text{(by partitioning } V_i \text{ into } C_1, \dots, C_{n_{i+1}} \text{ as the algorithm continues until it marks all vertices)} \\ &\geq n_{i+1} \cdot t_i. && \text{(by Eq (1))} \end{aligned}$$

Thus, we have that

$$\log t_{i+1} = \frac{2m_i}{n_{i+1}} \geq \frac{2m_i}{n_{i+1}} \geq t_i,$$

and hence $t_{i+1} \geq 2^{t_i}$, concluding the proof. \square

By Lemma 6, in round i of the algorithm, the value of t_i is at least a *tower* of twos of length i (i.e., two to the power two to the power two ... for i times). This means that after $O(\log^*(n))$ rounds, the value of t_i will be at least n and in that round, we are only running Prim's algorithm once and so the algorithm terminates surely in that round. Thus, we conclude Algorithm 4 takes $O(m \log^*(n))$ time.

This concludes the proof of Theorem 5.

References

- [FT87] Michael L. Fredman and Robert Endre Tarjan. Fibonacci heaps and their uses in improved network optimization algorithms. *J. ACM*, 34(3):596–615, 1987. 6, 7