

## Lecture 7

September 27, 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

<b>1</b>	<b>Karger-Klein-Tarjan Algorithm for MST</b>	<b>1</b>
1.1	Preliminaries . . . . .	1
1.2	The <i>KKT Algorithm</i> . . . . .	2

## 1 Karger-Klein-Tarjan Algorithm for MST

In the previous lecture, we went over the basics of the MSTs, the classical algorithms for it, and then saw the Fredman-Tarjan algorithm that solves this problem deterministically in  $O(m \log^*(n))$  time. While for all practical purposes, this is as good as a linear time algorithm, mathematically speaking,  $\log^*(n)$  still goes to  $+\infty$  with  $n$ , no matter how slowly, and thus, this runtime is still truly  $\omega(m)$ .

In this lecture, we will go over another algorithm for this problem due to Karger, Klein, and Tarjan [KKT95], which runs in  $O(m)$  time but it is randomized, i.e., its  $O(m)$  runtime is in expectation (and with high probability).

**Theorem 1** ([KKT95]). *There is a randomized algorithm for the minimum spanning tree problem that runs in  $O(m)$  time in expectation and with high probability.*

In this lecture, we only prove the runtime of the algorithm in expectation (although extending it to a with high probability bound is quite simple also). Before moving on, we recall the following two key rules in the design of MST algorithms:

**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.*

**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.*

The general idea behind the Karger-Klein-Tarjan algorithm (henceforth, *KKT algorithm*) is to use the **cycle rule** to get rid of most edges of the graph quickly, namely, *sparsify* the graph, and then solve the problem recursively on this sparser graph. To present this algorithm, we need some preliminaries first.

### 1.1 Preliminaries

The following definition is key to the design of *KKT algorithm*.

**Definition 2.** Fix any graph  $G = (V, E)$  and any forest  $F$  that is a *subgraph* of  $G$ . We say that an edge  $e \in E \setminus F$  is  **$F$ -heavy** if adding  $e$  to  $F$  results in a cycle and  $e$  is the maximum weight edge of that cycle. We refer to any other edge as an  **$F$ -light** edge.

The following observation is a direct corollary of the **cycle rule** and the definition of  $F$ -heavy edges.

**Observation 3.** For any graph  $G$ , any forest  $F$  that is a subgraph of  $G$ , and any  $F$ -heavy edge  $e$ , the MST of  $G - e$  is the same as the MST of  $G$ .

How do we use **Observation 3** in the algorithm? Suppose first that  $F$  is the MST of  $G$ ; then *every* edge  $e \in G \setminus F$  is  $F$ -heavy and thus can be neglected when computing the MST of  $G$ . Obviously however, this is not helpful as we need to first compute the MST of  $G$ . But, what if we have some other forest  $F$  which is easier to compute than the MST? Then, **Observation 3** tells us that we can still neglect *every*  $F$ -heavy edges without any worry. Thus, in order to find the MST of  $G$ , we can first try to quickly find “some approximate” forest  $F$ , with the key property that *most* edges of the graph are  $F$ -heavy, and then recursively solve the problem on the remaining few  $F$ -light edges. This is precisely what *KKT algorithm* does.

There is one more missing ingredient in the above approach. Given a graph  $G$  and a forest  $F$ , how quickly can we find the set of  $F$ -heavy edges? It is easy to check if a *single* each is  $F$ -heavy or not in linear time. But, doing this for every edge this way separately leads to a quadratic time algorithm which is way above our budget. Nevertheless, a surprising fact is the we can find *all*  $F$ -heavy edges in linear time as well!

**Theorem 4** ([Kom85, DRT92, Kin97]). There is an algorithm that given any graph  $G$  and any forest  $F$  which is a subgraph of  $G$ , outputs the set of all  $F$ -heavy edges in  $O(m + n)$  time.

We will not cover this algorithm in this course and just take it for granted. We only mention that there is a great deal of algorithmic work on this problem with the goal of finding simpler algorithms, but it does look like we still do not have a particularly simple algorithm for this problem. Also a note that these algorithms are often studied for the purpose of *MST verification* since having such an algorithm allows us to detect if a given tree  $F$  is an MST or not (the MST of the graph is the only one that makes all other edges  $F$ -heavy).

## 1.2 The KKT Algorithm

We are now ready to present the *KKT algorithm*. We note that given the recursive nature of the algorithm and since it can be called on not-necessarily connected graph, we use the term *Minimum Spanning Forest (MSF)* throughout which refers to a collection of MSTs on each connected components of the graph.

**Algorithm 1 (Karger-Klein-Tarjan Algorithm).**

- (i) Run 3 rounds of the Boruvka’s algorithm and let  $G'$  be the contracted graph obtained from  $G^a$ .
- (ii) Sample each edge of  $G'$  independently with probability  $1/2$  to obtain a graph  $G_1$ . **Recursively** find the MSF of  $G_1$  and call it  $F$ .
- (iii) Use the algorithm of **Theorem 4** to find all  $F$ -heavy edges of  $G'$  and let  $G_2$  be the graph obtained from  $G'$  after removing them.
- (iv) **Recursively** find the MSF of  $G_2$  and return it as the answer.

<sup>a</sup>This is a simple preprocessing step to reduce the number of vertices slightly

**Proof of Correctness.** The correctness of this algorithm is actually quite easy to proof. The first step is correct due to the correctness of Boruvka’s algorithm established earlier. Regardless of the choice of  $G_1$  and

the resulting MSF  $F$ , we have by [Observation 3](#) that none of the edges removed from  $G$  to obtain  $G_2$  can be part of the MSF of  $G$ . Thus, finding the MSF of  $G_2$  is the same as the MSF of  $G$  to begin with, and thus the algorithm returns the correct answer.

**Runtime Analysis.** The key to the analysis of the algorithm is to show that the graph  $G_2$  actually has few edges. In other words, after picking the MSF  $F$  on (almost) half the edges, the set of  $F$ -heavy edges more or less contains all but  $O(n)$  edges of the graph.

**Lemma 5.** *The expected number of  $F$ -light edges in [Algorithm 1](#) is at most  $2 \cdot (n' - 1)$  where  $n'$  is the number of vertices in  $G'$ .*

*Proof.* Notice that even though we are computing MSF of  $G_1$  using a recursive call to the *KKT algorithm*, given that MSF is unique (recall our assumption on distinct weights from the last lecture), for the purpose of the analysis, we can assume  $F$  is instead computed using Kruskal's algorithm. This is because the distribution of  $F$  is identical in both cases.

Now, let us examine how Kruskal's algorithm works. Suppose we sort all edges of  $G'$  (and not only  $G_1$ ) in increasing order of weight and call them  $e_1, \dots, e_{m'}$ . Consider the following process. We go over these edges one by one call  $F_i$  the subgraph of  $F$  maintained so far when visiting the edge  $e_i$ . We check if adding  $e_i$  to  $F_i$  creates a cycle or not. If it does, then whether or not  $e_i$  is sampled in  $G_1$  we are not going to pick this edge in the MSF  $F$  so we just ignore it. But, if it does not, it is only now that we check whether  $e_i$  belongs to  $G_1$  even or not. This means that only now we toss the coin to decide if  $e_i$  joins  $G_1$  or not. Notice that, despite all these seeming changes, we actually have not changed the distribution of  $F$  in anyway in this process (we can toss a coin for neglected edges and include them in  $G_1$  if we want just to make sure the distribution of  $G_1$  remains identical, although this does not change the distribution of  $F$  in any way).

Finally, note that all the edges ignored in this process are certainly  $F$ -heavy because they created a cycle even with a subgraph of  $F$  and are the heaviest weight edge of that cycle. Thus, the number of  $F$ -light edges is at most equal to the number of edges that we did not ignore, in other words, the edges that we tossed a coin for. At the same time, whenever we toss a coin, with probability half, we add the edge to the forest  $F$ . Moreover, the forest  $F$  cannot have more than  $n' - 1$  edges. So, the expected number of coin tosses we can have before collecting  $n' - 1$  edges in  $F$  is  $2 \cdot (n' - 1)$ , proving the lemma.  $\square$

We are now ready to conclude the proof. Firstly, let  $A(G, r)$  denote the runtime of the algorithm on a graph  $G$  when *all* the random bits we use is  $r$  (note that  $A(G, r)$  is deterministically fixed after we fixed the randomness). We have,

$$A(G, r) \leq c \cdot (m + n) + A(G_1, r) + A(G_2, r),$$

for some absolute constant  $c > 0$  which is the hidden constant in  $O(m + n)$  time needed for Boruvka's algorithm in the first step, the use of [Theorem 4](#), and general bookkeeping throughout the algorithm ignoring the recursive calls. Thus, the expected runtime of the algorithm on a graph  $G$  is

$$\mathbb{E}_r[A(G, r)] \leq c \cdot (m + n) + \mathbb{E}_r[A(G_1, r)] + \mathbb{E}_r[A(G_2, r)]. \quad (1)$$

Now, define  $T(m, n, r)$  as the worst-case runtime of the algorithm on a graph with  $m$  edges,  $n$  vertices, and for the randomness  $r$ . We prove inductively that

$$\mathbb{E}_r[T(m, n, r)] \leq 2c \cdot (m + n).$$

For any graph  $H$ , let  $m(H)$  and  $n(H)$ , denote the number of edges and vertices in  $H$ , respectively. Also, let  $r_1$  be the randomness used out of the recursive calls and  $r_2$  be the randomness of the recursive calls. Given [Eq \(1\)](#), we have,

$$\mathbb{E}_r[T(m, n, r)] \leq c \cdot (m + n) + \mathbb{E}_r[T(m(G_1), n(G_1), r)] + \mathbb{E}_r[T(m(G_2), n(G_2), r)]$$

$$\begin{aligned}
&\leq c \cdot (m + n) + \mathbb{E}_{r_1} \left[ \mathbb{E}_{r_2} [T(m(G_1), n(G_1), r_2)] \right] + \mathbb{E}_{r_1} \left[ \mathbb{E}_{r_2} [T(m(G_2), n(G_2), r_2)] \right] \\
&\quad \text{(the recursive calls themselves only depend on } r_2 \text{ (after fixing their input based on } r_1)) \\
&\leq c \cdot (m + n) + \mathbb{E}_{r_1} [2c \cdot (m(G_1) + n/8)] + \mathbb{E}_{r_1} [2c \cdot (m(G_2) + n/8)] \\
&\quad \text{(by induction hypothesis and as 3 rounds of Boruvka's algorithm reduces vertices by a factor of 8)} \\
&\leq c \cdot (m + n) + 2c \cdot (m/2 + n/8) + 2c \cdot (2n/8 + n/8) \\
&\quad \text{(as } \mathbb{E}[m(G_1)] = m'/2 \leq m/2 \text{ trivially and } \mathbb{E}[m(G_2)] \leq 2n' \leq 2n/8 \text{ by Lemma 5)} \\
&= 2c \cdot m + c \cdot (n + n/4 + 3n/4) = 2c \cdot (m + n),
\end{aligned}$$

proving the induction step. Thus, the runtime of the algorithm is  $O(m + n)$  in expectation.

This concludes the proof of [Theorem 1](#) and our study of MST algorithms in this course.

**Remark.** The *KKT algorithm* provided the first linear time algorithm for MSTs but at the “cost” of randomization. Hence, the search for a deterministic algorithm for this problem still continues and to date we do not know such an algorithm. The Fredman-Tarjan algorithm we discussed in the last lecture has been improved multiple times, with the current best deterministic algorithm due to Chazelle [[Cha00](#)] with runtime  $O(m \cdot \alpha(n))$  where  $\alpha(n)$  is a certain *Inverse Ackerman* function. There is also the algorithm of Pettie and Ramachandran [[PR02](#)] that is *provably optimal* but its runtime is not known.

## References

- [Cha00] Bernard Chazelle. A minimum spanning tree algorithm with inverse-ackermann type complexity. *J. ACM*, 47(6):1028–1047, 2000. [4](#)
- [DRT92] Brandon Dixon, Monika Rauch, and Robert Endre Tarjan. Verification and sensitivity analysis of minimum spanning trees in linear time. *SIAM J. Comput.*, 21(6):1184–1192, 1992. [2](#)
- [Kin97] Valerie King. A simpler minimum spanning tree verification algorithm. *Algorithmica*, 18(2):263–270, 1997. [2](#)
- [KKT95] David R. Karger, Philip N. Klein, and Robert Endre Tarjan. A randomized linear-time algorithm to find minimum spanning trees. *J. ACM*, 42(2):321–328, 1995. [1](#)
- [Kom85] János Komlós. Linear verification for spanning trees. *Comb.*, 5(1):57–65, 1985. [2](#)
- [PR02] Seth Pettie and Vijaya Ramachandran. An optimal minimum spanning tree algorithm. *J. ACM*, 49(1):16–34, 2002. [4](#)