

# CS31 (Algorithms), Spring 2020 : Lecture 10

Date:

Topic: Graph Algorithms 1: Depth First Search

*Disclaimer: These notes have not gone through scrutiny and in all probability contain errors.*

*Please discuss in Piazza/email errors to deeparnab@dartmouth.edu*

---

## 1 Depth First Search (DFS)

We start graph algorithms with the pretty intuitive, but surprisingly powerful, *depth first search* (DFS). This algorithm solves the reachability problem, but then in one swoop solves much more. It also runs in  $O(n+m)$  time. Let's get to it.

### 1.1 DFS from a vertex

Our journey starts from a given vertex  $v$  in the graph. As explorers in the past have done, we begin by first planting a “flag” at  $v$ . In algorithmic terms, we have a Boolean variable  $\text{visited}[x]$  for every vertex  $x$  initialized to 0 (false). When starting at  $v$ , our first action is to set  $\text{visited}[v] = 1$ . Subsequently, we investigate  $v$ 's out-neighbors, and if we encounter any “unflagged” neighbor  $x$ , we proceed to  $x$  remembering (via a thread, say, to give a physical mnemonic) that we came from  $v$ , and then just repeat the procedure from  $x$ . If at some point we see that all of  $x$ 's neighbors are flagged, then we ravel the metaphorical thread back to the place where  $x$  was called from (in this case  $v$ ), and then continue the investigation of  $v$ 's other neighbors. All this thread business is sweetly handled by recursion. Here is the pseudocode for DFS from a vertex.

```
1: global visited[1 : n] initialized to 0.
2: global Tree  $F$  initialized to  $\perp$ .
3: procedure DFS( $G, v$ ):  $\triangleright$  We assume  $V = \{1, 2, \dots, n\}$ 
4:    $\triangleright$  Returns a tree  $F$  rooted at  $v$ 
5:   visited[ $v$ ]  $\leftarrow$  1.  $\triangleright$  Mark  $v$  visited.
6:   Add vertex  $v$  to  $F$ .
7:   for  $u$  neighbor of  $v$  do:  $\triangleright$  In an arbitrary, but fixed order
8:     if visited[ $u$ ] = 0 then:
9:       Add edge  $(v, u)$  to the tree  $F$ .
10:      DFS( $G, u$ )
11:   return  $F$ 
```

As you can see from the pseudocode, there is some extra stuff that the algorithm is building. Namely, the *global* object  $F$ . It begins with being  $\perp$ , and then in [Line 9](#) a bunch of edges are added to  $F$ . These edges are indeed the “threads” that were alluded to above. At the end of the call of  $\text{DFS}(G, v)$ , the algorithm returns this collection of edges. The following theorem shows that this collection indeed has a lot of structure, and contains information about the *connectivity* properties of  $G$ .

**Theorem 1.** A vertex  $x \in F$  if and only if  $x$  is reachable from  $v$  in  $G$ .  $F$  is a *rooted out-tree* rooted at  $v$ .

*Proof.* First we prove if  $\text{visited}[x] = 1$ , then  $x$  is reachable from  $v$ . In fact we will show that  $x$  is reachable from  $v$  in  $F$ . The proof is by induction on the *time* at which  $\text{visited}[x]$  was set to 1. Imagine every time the algorithm runs [Line 5](#), we increment time by 1. At time 0, this was set  $\text{visited}[v] = 1$  and  $v$  is reachable from  $v$  in  $F$ . Now pick a vertex  $x$  whose  $\text{visited}[x] = 1$  is set at time  $t$ . This happens because of some  $y \in V$  such that (a)  $(y, x) \in E$ , and (b) the run of  $\text{DFS}(G, y)$  calls  $\text{DFS}(G, x)$ . In that case, (a)  $\text{visited}[y] = 1$  has been set strictly before time  $t$  implying, by induction,  $y$  is reachable from  $v$  in  $F$ , and (b) we add edge  $(y, x)$  to  $F$ , which then implies  $x$  is reachable from  $v$  in  $F$ .

Now for the other direction. Suppose there exists a vertex  $x$  which is reachable from  $v$  in  $G$  but  $\text{visited}[x] = 0$ . Since  $x$  is reachable from  $v$ , there is a path  $(v = v_0, v_1, \dots, v_k = x)$  in  $G$ . Let us pick the *last* vertex  $v_i$  in this path which has  $\text{visited}[v_i] = 1$ ; clearly  $0 \leq i < k$ . Since  $\text{visited}[v_i] = 1$ , we have run  $\text{DFS}(G, v_i)$ . But the for-loop in the algorithm would then call  $\text{DFS}(G, v_{i+1})$  since  $\text{visited}[v_{i+1}] = 0$ . But that would set  $\text{visited}[v_{i+1}] = 1$ , and once visited a vertex is never “un-visited”. This is a contradiction, and thus all vertices  $x$  reachable from  $v$  have  $\text{visited}[x] = 1$ .

The first part above shows that for every  $x \in F$ , there is a path from  $v$  to  $x$  in  $F$ . To show  $F$  is a rooted out-tree, we need to show that if we ignore the directions of the edges of  $F$ , the underlying undirected graph is a tree. Suppose not, and suppose there is a cycle  $C$  in the underlying undirected graph of  $F$ . Let the edge  $(x, y)$  be the *last* directed edge being added to this cycle. That is,  $\text{DFS}(G, x)$  called  $\text{DFS}(G, y)$ . In particular, at this time,  $\text{visited}[y] = 0$  (otherwise, the algorithm wouldn’t call  $\text{DFS}(G, y)$ ). However, there is an edge  $(y, z)$  in the (undirected) cycle  $C$ . That is, either  $(y, z) \in F$  or  $(z, y) \in F$ . And this has occurred before  $\text{DFS}(G, x)$  is going to call  $\text{DFS}(G, y)$ . However, the presence of this edge (either  $(y, z)$  or  $(z, y)$ ) *implies*  $\text{visited}[y] = 1$ . This is a contradiction since  $\text{visited}[y] = 0$ . This proves that  $F$  is a rooted out-tree rooted at  $v$ .  $\square$

**Theorem 2.** The REACHABLE? problem can be solved in  $O(n + m)$  time.

*Proof.* [Theorem 1](#) implies this as a corollary. Given vertex  $u$  and  $v$ , we can check if  $v$  is reachable from  $u$  by just running  $\text{DFS}(G, u)$ , and checking if  $\text{visited}[v] = 1$  or not. To get the path, we can use the tree  $F$ . The running time is  $O(n + m)$  since every edge is considered in the for-loop at most twice (once if  $G$  is directed).

What is the running time of DFS? It is a recursive algorithm, so it is not trivial to see this. Let us fix an arbitrary vertex  $x$ , and let us figure out the time taken in the running of its pseudocode *other* than the DFS calls it makes. That is, the time for running [Line 5](#), [Line 6](#), [Line 9](#). We see that the maximum time is taken in the for-loop, and this costs  $O(\text{deg}^+(x))$  time. The second, and the key, observation is that in *all* the calls of  $\text{DFS}(G, v)$ , a call  $\text{DFS}(G, x)$  for any vertex  $x$  is made at most once (exactly once for vertices reachable from  $v$ ). This is because, once  $\text{DFS}(G, x)$  is called,  $\text{visited}[x]$  is set to 1 which prevents any further calls. Thus, the total time taken by  $\text{DFS}(G, v)$  is at most  $\sum_{x \in V}$  the total time taken by “non-recursive” calls of  $\text{DFS}(G, x)$ , which is  $\sum_{x \in V} O(\text{deg}^+(x))$ . This evaluates to  $O(n + m)$ .  $\square$

## 1.2 DFS on the whole graph

The next algorithm is a *traversal* over all vertices of the graphs using the subroutine  $\text{DFS}(G, v)$  repeatedly. This is called the *depth first traversal* algorithm of the graph  $G$ , but is also called the depth first search (or simply DFS) of  $G$ . The input to this algorithm is the graph  $G$  and a permutation/ordering  $\sigma$  of the vertices. This permutation tells the algorithm the order in which to “explore” vertices, that is, to run  $\text{DFS}(G, v)$ .

The output to this algorithm has a lot of things; these objects contain surprising amounts of information about  $G$ , as we will see below.

- One output is a couple of vectors  $\text{first}[1 : n]$  and  $\text{last}[1 : n]$  where for any vertex  $v$ ,  $\text{first}[v]$  notes the “time” at which the algorithm starts exploring from  $v$ , that is,  $\text{DFS}(G, v)$  is called, and  $\text{last}[v]$  denotes the “time” the exploring ends, that is, the subroutine  $\text{DFS}(G, v)$  terminates.
- The other output is a *forest*  $F$  spanning all the vertices of  $G$ . Each tree in the forest is a rooted out-tree, and is assumed to be directed (even when  $G$  is not) away from the root. Together with this we store the scalar  $\text{fcomp}$  which counts the number of trees in  $F$ , the array  $\text{root}[1 : \text{fcomp}]$  where  $\text{root}(i)$  will store the root of the  $i$ th tree in  $F$ , and the array  $\text{Fcomp}[1 : n]$  where  $\text{Fcomp}[v]$  contains a number between 1 and  $\text{fcomp}$  indicating the tree in which  $v$  exists.

The algorithm is simple: it has a for-loop going over all vertices in the order  $\sigma$ ; if the vertex is unvisited, then we run  $\text{DFS}(G, v)$  on it starting a new tree rooted from  $v$ . We end when there are no more vertices left. Here is the full pseudocode, where we have enhanced  $\text{DFS}(G, v)$  to take care of what we need.

```

1: procedure  $\text{DFS}(G, \sigma[1 : n])$ :  $\triangleright \sigma$  is an ordering of the vertices
2:   global array  $\text{visited}[1 : n]$  initialized to all 0.
3:   global array  $\text{first}[1 : n], \text{last}[1 : n], \text{root}[1 : n], \text{Fcomp}[1 : n]$  initialized to all 0.
4:   global scalar  $\text{fcomp}$ , time initialized to 0.
5:   global Forest  $F$  initialized to  $\emptyset$ .

6:   for  $v$  in  $\sigma$  do:
7:     if  $\text{visited}[v] = 0$  then:  $\triangleright v$  hasn't been visited yet:
8:        $\text{fcomp} \leftarrow \text{fcomp} + 1$   $\triangleright$  Increase the number of trees in the forest
9:        $\text{root}[\text{fcomp}] \leftarrow v$   $\triangleright$  Set  $v$  to be the root of the new tree
10:       $\text{DFS}(G, v)$ 

11: procedure  $\text{DFS}(G, v)$ :
12:    $\text{visited}[v] \leftarrow 1$ ; Add  $v$  to  $F$ .
13:    $\text{Fcomp}[v] \leftarrow \text{fcomp}$ .  $\triangleright$  Set  $v$ 's tree in the forest
14:    $\text{time} \leftarrow \text{time} + 1$ .
15:   Set  $\text{first}[v] \leftarrow \text{time}$ .  $\triangleright$  Start exploring.
16:   for  $u$  neighbor of  $v$  do:  $\triangleright$  In an arbitrary, but fixed order
17:     if  $\text{visited}[u] = 0$  then:
18:       Add edge  $(v, u)$  to the forest  $F$ .
19:        $\triangleright$  It will be added to the  $\text{fcomp}$ th component.
20:        $\text{DFS}(G, u)$ 
21:    $\text{time} \leftarrow \text{time} + 1$ .
22:   Set  $\text{last}[v] \leftarrow \text{time}$ .

```

As we argued in the case of  $\text{DFS}(G, v)$ , one can see that  $\text{DFS}(G, \sigma)$  takes  $O(n + m)$  time as well.

**Claim 1.** The running time of  $\text{DFS}(G, \sigma)$  is  $O(m + n)$  for any  $\sigma$ .

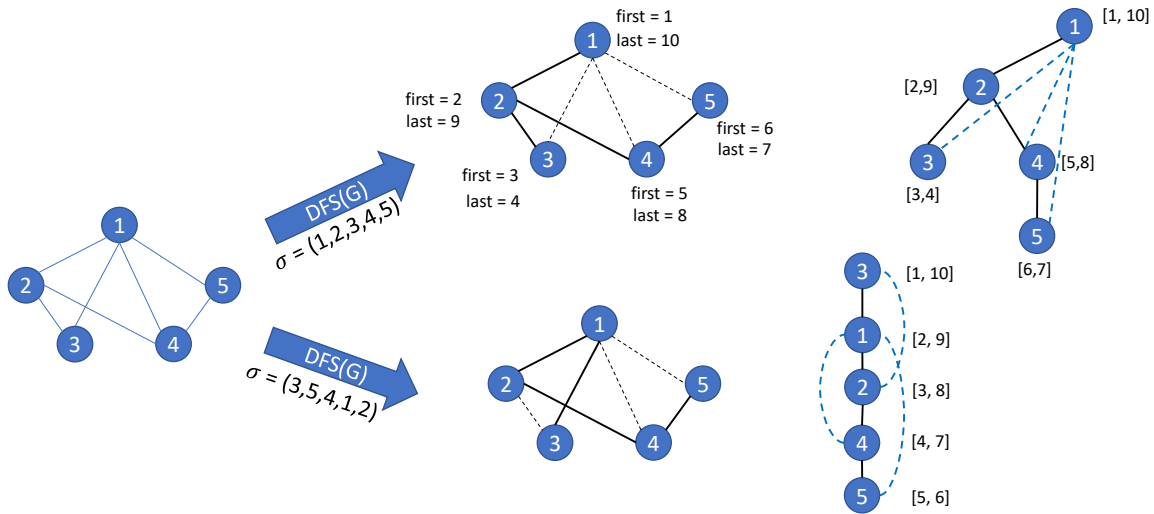


Figure 1: The edges that appear in the forest are marked in solid, while the remaining edges are dotted. The first and last are noted near the vertices. In the third figure on the right, the interval is the  $[first[v], last[v]]$  interval.

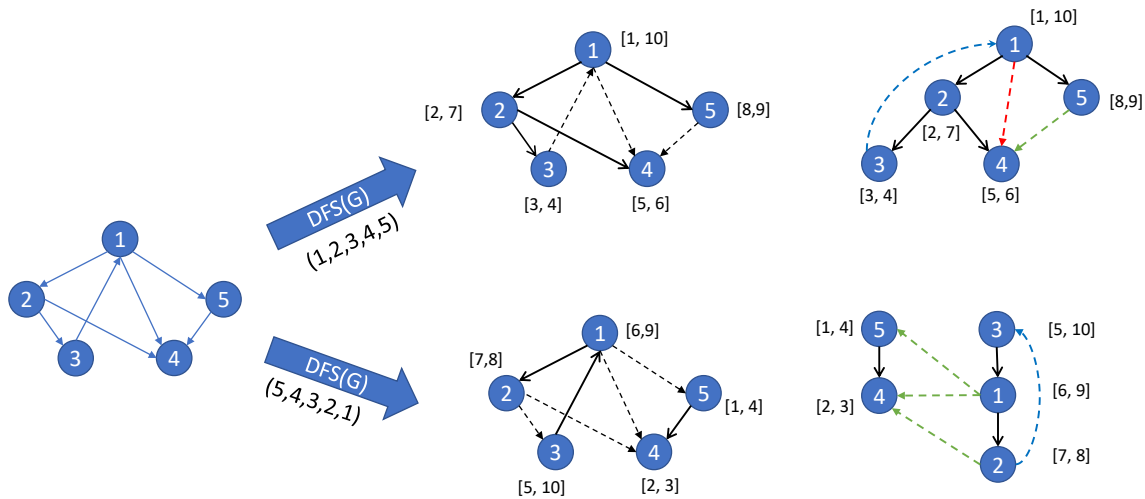


Figure 2: The edges that appear in the forest are marked in solid, while the remaining edges are dotted. The first and last are noted near the vertices. In the third figure on the right, the interval is the  $[first[v], last[v]]$  interval.

**Remark:** Different permutations can lead to different outcomes: see [Figure 1](#) and [Figure 2](#). This will be *critically* used in one application of DFS.

**Edge Classification.** Running DFS on a graph  $G$  with any ordering  $\sigma$  leads to four kinds of edges.

- (Forest Edges.) These are the edges present in  $F$ . These are marked black and solid in the Figures.
- (Back Edges.) These edges go from a descendant to an ancestor. These are marked blue and dotted.
- (Forward Edges.) These edges go from an ancestor to a descendant. These are marked red and dotted. For undirected graphs the forward edges are all back edges (there is no direction).
- (Cross Edges.) All the rest. They can be among pairs in the same component, or not. These are marked green and dotted.

**Properties.** Next, we state and prove **three** properties of the output of the DFS algorithm. Before reading the proofs, it will be useful to see their illustrations in the examples shown in [Figure 1](#) and [Figure 2](#) (or any other figures you have privately made). For all the properties below, we assume we have run  $\text{DFS}(G, \sigma)$  for some *arbitrary* ordering  $\sigma$ .

**Lemma 1 (Nested Interval Property).** For any two vertices  $u$  and  $v$ , with  $\text{first}[u] < \text{first}[v]$ , exactly one of the following two properties hold.

- $\text{first}[u] < \text{first}[v] < \text{last}[v] < \text{last}[u]$  and  $v$  is a descendant of  $u$  in  $F$ .
- $\text{first}[u] < \text{last}[u] < \text{first}[v] < \text{last}[v]$  and neither is a descendant of the other.

This shows that the  $n$  intervals of the form  $[\text{first}[v], \text{last}[v]]$  don't "criss-cross" (although one may be contained in the other). This property is called the *nested* property (also called laminar property).

*Proof.* Since  $\text{first}[u] < \text{first}[v]$ , we call  $\text{DFS}(G, u)$  before we call  $\text{DFS}(G, v)$ . If  $v$  is ever discovered in the run of  $\text{DFS}(G, u)$ , then (a) it will be a descendant of  $u$  in  $F$  (by [Theorem 1](#)), and (b) we must have  $\text{last}[v] < \text{last}[u]$  since this recursive call must end before  $u$ 's recursive call ends. This is case 1.

The only other case is  $v$  has not been discovered in the run of  $\text{DFS}(G, u)$ . In that case,  $\text{last}[u] < \text{first}[v]$  by definition. Furthermore,  $v$  is not a descendant of  $u$ , and  $u$  can't be a descendant of  $v$  since  $v$  hasn't been even discovered yet. This is case 2. □

The above property is useful, and will be useful in proving some other properties below. But it also allows us to classify the edges (not in the forest  $F$ ) just looking at the first and the last values.

- (Back Edges.) Edges  $(u, v) \in E \setminus F$  with  $\text{first}[v] < \text{first}[u] < \text{last}[u] < \text{last}[v]$ .
- (Forward Edges.) Edges  $(u, v) \in E \setminus F$  with  $\text{first}[u] < \text{first}[v] < \text{last}[v] < \text{last}[u]$ .
- (Cross Edges.) Edges  $(u, v) \in E$  such that the intervals  $[\text{first}[u], \text{last}[u]]$  and  $[\text{first}[v], \text{last}[v]]$  are disjoint.

**Lemma 2 (Edge Property).** Let  $(u, v)$  be any edge in  $G$  with  $\text{first}[u] < \text{first}[v]$ . Then, we must have  $\text{last}[v] < \text{last}[u]$ .

*Proof.* Suppose not. By the Nested Interval Property, we must have  $\text{first}[u] < \text{last}[u] < \text{first}[v] < \text{last}[v]$ . That happens when  $\text{DFS}(G, u)$  terminates before  $\text{visited}[v]$  is set to 1. But the [Line 16](#) would discover  $v$  contradicting the above. □

We are now ready for the first application of DFS – we can solve CONNECTED COMPONENTS of an Undirected Graph using the following lemma.

**Lemma 3.** Let  $G = (V, E)$  be any undirected graph and consider the forest  $F$  returned by  $\text{DFS}(G, \sigma)$  with any permutation  $\sigma$ . The components of  $F$  are precisely the connected components of  $G$ .

*Proof.* Let  $V_1, \dots, V_k$  be the vertices in the various trees of the forest  $F$ . Clearly  $G[V_i]$  is connected since they are connected in the forest. We claim that there is no edge of the form  $(u, v)$  with  $u \in V_i$  and  $v \in V_j$ . Suppose there is, and without loss of generality assume  $\text{first}[u] < \text{first}[v]$  (this is where we are using the undirectedness of  $G$ ). By the edge property, we have  $\text{first}[u] < \text{first}[v] < \text{last}[v] < \text{last}[u]$ . But this means  $v$  is a descendant of  $u$  in  $F$  contradicting the fact they exist in different connected components of  $F$ .  $\square$

**Theorem 3.** CONNECTED COMPONENTS of an undirected graph can be found in  $O(n + m)$  time by running  $\text{DFS}(G, \sigma)$  for any ordering  $\sigma$ .

Moving on to more properties.

**Lemma 4 (Path Property).** If  $(u = v_1, v_2, \dots, v_k = v)$  is a path in  $G$  from  $u$  to  $v$  such that  $\text{first}[u] < \text{first}[v_i]$  for all  $2 \leq i \leq k$ , then,  $\text{last}[v_i] < \text{last}[u]$  for all  $2 \leq i \leq k$ .

In English, if there is a path from a vertex  $u$  to a vertex  $v$  such that  $u$  is the first vertex to be discovered among them, then all the vertices in the path are descendants of  $u$  in the DFS forest.

*Proof.* Suppose not. Choose the smallest  $2 \leq i \leq k$  for which  $\text{last}[u] < \text{last}[v_i]$ . By the choice of  $i$ , we get  $\text{last}[v_{i-1}] < \text{last}[u]$ . Also note  $(v_{i-1}, v_i)$  is an edge.

Case 1:  $\text{first}[v_{i-1}] < \text{first}[v_i]$ . In this case, the Edge Property would imply  $\text{last}[v_i] < \text{last}[v_{i-1}]$ , and thus  $\text{last}[v_i] < \text{last}[u]$ . Which is what we supposed wasn't true.

Case 2:  $\text{first}[v_i] < \text{first}[v_{i-1}]$ . In that case, we see  $\text{first}[u] < \text{first}[v_i] < \text{last}[u] < \text{last}[v_i]$  which violates the Nested Interval Property.  $\square$

The above property allow us immediately to solve the CYCLE? problem. The following theorem implies the algorithm: run  $\text{DFS}(G, \sigma)$  and check if any of the edges is a *back edge* (which is one linear time scan over all the edges and checking the first and the last).

**Lemma 5.** A graph  $G$  is **acyclic** if and only if there are no back edges.

*Proof.* One direction is trivial – if  $G$  has a back edge, then there is clearly a cycle. If the back-edge is  $(u, v)$ , then by definition there is a path from  $v$  to  $u$  using  $F$ -edges, and then take the  $(u, v)$  edge back.

The other direction is more interesting. If  $G$  has a cycle  $C$  with  $k$  vertices  $(v_1, \dots, v_k, v_1)$ , then without loss of generality let  $v_1$  be the vertex with the smallest  $\text{first}[v_i]$  in this cycle. Since there is a path from  $v_1$  to  $v_k$ , using the Path property and the fact that  $\text{first}[v_1]$  is the smallest, we get  $\text{first}[v_1] < \text{first}[v_k] < \text{last}[v_k] < \text{last}[v_1]$ . But this implies  $(v_k, v_1)$  is a back-edge.  $\square$

**Theorem 4.** CYCLE? can be solved in  $O(n + m)$  time using DFS.