

CS31 (Algorithms), Spring 2020 : Lecture 11

Date:

Topic: Graph Algorithms 2: Applications of DFS

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 Applications of DFS

We already saw two applications of DFS in the last lecture: the REACHABLE? problem and the CYCLE? problem. In this lecture, we see two more applications of DFS on directed graphs. The first application actually gives a way to “order” all nodes in a *directed acyclic graph* (DAG). This is called the *topological order*. This itself has many applications, which your problem sets explore. The second application is to figuring out whether a directed graph is *strongly* connected or not.

1.1 Topological Ordering of DAGs

Throughout this subsection, G is a directed acyclic graph (DAG). Recall from the previous lecture, this means that if we run DFS on G , there are no back edges. A topological ordering of (the vertices of) a DAG is an ordering $\sigma[1 : n]$ of the vertices such that for any $i < j$, there is **no** edge from $\sigma[j]$ to $\sigma[i]$. That is, if we write down the vertices from left to right in the σ order, then *all* edges go from left to right. If one thinks of an edge (u, v) as v being “bigger” than u , then the topological ordering is a linearization of the graph according to this (partial) order. Of course, not every pair of vertices may be comparable.

Remark: Note that the first vertex v in the topological order must have $\deg^-(v) = 0$. There is no vertex to its left to “send” an edge to it. Such vertices are called sources. Similarly, the last vertex v in the topological order must have $\deg^+(v) = 0$. There is no vertex to its right to which it can “send” an edge. Such vertices are called sinks.

Does a topological order always exist of a DAG? Before moving on, think a bit about this question. (Hint: maybe look at the question of whether DAGs have sources and sinks. Then consider removing them...)

TOPOLOGICAL ORDERING

Input: Directed Acyclic Graph G .

Output: A topological ordering of G .

If you haven’t been able to puzzle out whether a topological order exists, then you will probably appreciate the power of DFS.

Lemma 1. Consider running DFS in any arbitrary order on the DAG G . Let σ be the ordering of the vertices in *decreasing* order of $\text{last}[v]$. Then, σ is a topological order.

Proof. To show σ is a topological order, we need to show there is no edge going from right to left. To this end, pick two vertices x and y such that $\sigma[x] < \sigma[y]$. We need to show (y, x) is *not* an edge. Suppose, for contradiction’s sake, it is. Now, $\sigma[x] < \sigma[y]$ implies, by our definition, $\text{last}[y] < \text{last}[x]$. The *edge property* now implies that $\text{first}[x] < \text{first}[y]$. Why? If not, that is, if $\text{first}[y] < \text{first}[x]$, then since (y, x) is an edge

the edge property would imply $\text{last}[x] < \text{last}[y]$. But the reverse is true. Therefore, $\text{first}[x] < \text{first}[y]$. In sum, we have $\text{first}[x] < \text{first}[y] < \text{last}[y] < \text{last}[x]$. The Nested Interval Property now tells us y must be a descendant of x , that is, (y, x) is a back-edge. This contradicts that G is acyclic. \square

Theorem 1. TOPOLOGICAL ORDERING of any DAG G can be found in $O(n + m)$ time.

Proof. There is one extra thing needed to argue about following the previous lemma. We need the (decreasing) sorted order of lasts; how do we get that in $O(n + m)$ time. Two answers.

One, the $\text{last}[v]$'s are integers between 1 and $2n$; so we can sort in $O(n)$ time using Count-Sort. Two, and this is less modular, but we can read out the numbers in increasing order of $\text{last}[v]$ as DFS as being run; after all, it is precisely the order in which the $\text{last}[v]$'s are set. The topological order is the reverse of that. \square

One application of the topological ordering is to decide whether there is a path in a DAG which contains *all* the vertices: simply take the topological ordering, and check if that is a path. If it is, we are done. Else, I claim there *cannot* exist such a path. If there were, then one of the edges of that path would be in the “wrong direction.” I leave the details for you to figure out.

1.2 Strongly Connected Components (SCCs) using DFS

We now see that the *strongly connected components* of a graph G can be found in *linear* $O(n + m)$ time. This is a truly surprising algorithm which really illustrates the power of DFS. Let's try to first sketch the main ideas behind the algorithm, and then subsequently give the final description and analysis. To do so, consider the graph in [Figure 1](#).

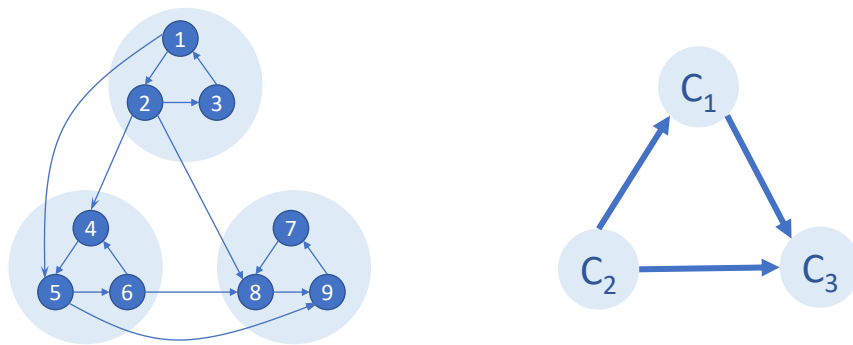


Figure 1: Illustrative example for SCC.

In the graph on the left, there are three strongly connected components marked in light blue circles. The graph on the right is one whose vertices are these three components, and we have an edge between two components (for instance from C_1 to C_2) if and only if there is an edge (u, v) in the original graph with $u \in C_1$ and $v \in C_2$. Note there could be, and in this example there are, multiple such edges. We require that there be at least one. In general, given a graph G as on the left, then the graph on the right is called G^{SCC} ; note we don't have this graph up front but is useful for analysis and designing the algorithm.

Claim 1. G^{SCC} is a DAG.

Proof. Suppose not, and suppose there was a cycle $(C_1, C_2, \dots, C_k, C_1)$ in G^{SCC} . This means there are vertices x_i, y_i in C_i (possibly x_i and y_i are the same) such that (x_i, y_{i+1}) is an edge in G , and (x_k, y_1) is also an edge in G . But then, we argue next, $C_1 \cup \dots \cup C_k$ should have been one strongly connected component, which would be a contradiction. Take any two vertices u and v in the union. Say $u \in C_i$ and $v \in C_j$ where $i \leq j$ without loss of generality (otherwise, we swap names). We now show a path from u to v . First, we go from u to x_i which is possible since C_i is strongly connected, then we take the (x_i, y_{i+1}) edge to y_{i+1} , and from there to x_{i+1} (since C_{i+1} is strongly connected), and so on and so forth till we reach y_j , upon reaching which, we take the path from y_j to v (which again exists since C_j is strongly connected). \square

Before we move on to discovering the SCCs, let us see why the algorithm for undirected graphs is not enough. Recall what we did for undirected graphs; we ran DFS on G in *any arbitrary* order and returned the connected components of the forest. Why doesn't it work? Well, in the graph in [Figure 1](#) consider what happens when we run DFS from the vertex 1. You see that *all* the vertices are reachable from 1 and thus end up in the tree rooted at 1. The resulting vertices are not strongly connected. To stress why this is not an issue in undirected graphs note that in undirected graphs if there is a path from a vertex 1 to a set of vertices S , then there is a path from any vertex in S to 1 as well. This is patently false in directed graphs.

An Encouraging Idea. Suppose that in the graph in [Figure 1](#), we ran DFS from vertex number 9. Then, we would definitely discover all the vertices that 9 can reach. But these are precisely the ones in C_3 , the strongly connected component connecting 9. Why is this? This is because, there is no edge which starts from inside C_3 and goes outside. That is, because C_3 is a *sink* component of G^{SCC} . But this is wonderful; there is some vertex from which if we start DFS we get at least one strongly connected component. Let us make this (finding one strongly connected component) our goal for now.

From our understanding of topological ordering in DAGs, we know that the vertex with the smallest $\text{last}[v]$ is a sink vertex in a DAG. Perhaps, we could conjecture something similar for a general graph: in any graph G and any DFS run, the vertex with the smallest $\text{last}[v]$ must lie in a sink component of G^{SCC} . If that is the case, then we could get what we want. *Unfortunately, this is not true.* Consider the graph in [Figure 1](#) again, with DFS being run in the $\{1, 2, \dots, 9\}$ order, and the adjacency lists also being visited in this order. We see that vertex 3 has the smallest $\text{last}[\]$, and indeed 3 lies in a source component of G^{SCC} .

Remark: *A philosophical interlude. In research, we often think we have a good understanding of objects, and this leads us to make some conjectures. Just like we did above. And often they are wrong. I'll not lie – disappointment is usually the first response. But what really defines a researcher is resilience. Counterexamples are the world's ways of telling us, "Your understanding was incomplete. Refine them. Think harder." And when we do get back to the drawing board, or square one, the world often rewards us with epiphanies.*

Epiphany 1. Although the vertex with the *smallest* $\text{last}[v]$ may not be in a SINK component of G^{SCC} , it is in fact true that the vertex with the *largest* $\text{last}[v]$ *does indeed* lie in the SOURCE component of G^{SCC} . Again going back to the example in G^{SCC} , we see that the vertex 1 has the *largest* last , and it is in the source component C_1 of G^{SCC} .

In fact, more is true. For any component $C \in G^{\text{SCC}}$, define

$$f(C) = \max_{v \in C} \text{last}[v]$$

That is, $f(C)$ is the largest last in that component.

Lemma 2. If (C_i, C_j) is an edge in G^{SCC} , then $f(C_i) > f(C_j)$.

Before we prove the above lemma, let us see why it implies that the vertex with the largest last must lie in a source component of G^{SCC} . Suppose x is the vertex with the largest last, and suppose it lies in component C_j . Clearly, $f(C_j) = \text{last}[x]$, for x is the largest last vertex. Now if C_j were not a source component, there would be some component C_i with (C_i, C_j) an edge in G^{SCC} . The above lemma would imply $f(C_i) > f(C_j) = \text{last}[x]$. That is, there is a vertex $y \in C_i$ with $\text{last}[y] > \text{last}[x]$. That contradicts the choice of x . Thus, x must lie in a source component. Let us now prove the lemma.

Proof. First we make a claim about strongly connected components.

Claim 2. For any strongly connected component C if $x \in C$ has the largest last, then it also has the smallest first. In particular, x 's interval contains all the intervals of every other vertex in C .

Proof. Suppose not. Suppose $y \in C$, $y \neq x$ has the smallest first. Since C is strongly connected, there is a path from y to x . y has the smallest first among all vertices in this path, so by the *path property*, it must have the the largest last among all vertices in this path. In particular, $\text{last}[y] > \text{last}[x]$. Contradiction. \square

The proof of the lemma is as follows. Let x be the vertex in C_i with the largest last and y be the vertex in C_j with the largest last. Thus, $f(C_i) = \text{last}[x]$ and $f(C_j) = \text{last}[y]$. For the sake of contradiction, suppose $\text{last}[x] < \text{last}[y]$. By the Nested Interval Property, either (a) $\text{first}[y] < \text{first}[x]$, that is, x 's interval is completely contained in y 's interval, or (b) $\text{last}[x] < \text{first}[y]$, that is x 's interval is disjoint and lies before y 's interval.

We will reach a contradiction in both cases. Case (a) is easy: if x 's interval is completely contained in y 's interval, then there is a path from y to x in the DFS forest. In particular, that would imply an edge from C_j to C_i in the G^{SCC} contradicting the DAG nature of G^{SCC} .

In Case (b), x 's interval finishes before y 's interval. By the claim above, this means that the interval of *every* vertex in C_i finishes before the interval of *any* vertex in C_j starts. Now since (C_i, C_j) is an edge, there is some $u \in C_i$ and $v \in C_j$ such that (u, v) is an edge in G . From the claim, we see $\text{first}[u] < \text{last}[u] < \text{first}[v] < \text{last}[v]$; this contradicts the *edge property*. \square

So, the above lemma gives us a way to recognize a vertex in the *source* component of G^{SCC} . But we needed a vertex in the *sink* component of G^{SCC} . How are we going to get that? A second epiphany answers this.

Epiphany 2. Let G_{rev} be the graph where all edges of G have been reversed. Observe that the strongly connected components of G_{rev} are the precisely the same as those in G , and that $(G_{\text{rev}})^{\text{SCC}} = (G^{\text{SCC}})_{\text{rev}}$. In other words, the source components of $(G)^{\text{SCC}}$ are precisely the sink components of $(G_{\text{rev}})^{\text{SCC}}$. Therefore, if we run DFS on G and look at the vertex with the largest last $\text{last}[v]$ that is guaranteed to be in the *sink* component of $(G_{\text{rev}})^{\text{SCC}}$. Which will allow us to find the strongly connected components of (G_{rev}) . Which is the same as the strongly connected components of G . Done!

Strongly Connected Component Algorithm. We now state the full algorithm. The algorithm takes input G and returns the components of G^{SCC} . Furthermore, they are returned in the *topological order* in G^{SCC} . This fact is also useful (see a UGP problem regarding this).

- 1: **procedure** STRONGCONNCOMP(G):
- 2: ▷ *Returns the strongly connected components of G*
- 3: ▷ *The order in which it is returned is the Top. Ord. of G^{scc} .*
- 4: Run DFS($G, \{1, 2, \dots, n\}$) to get last[v] for every vertex.
- 5: π be the decreasing order of last[v]’s. ▷ *Can be found in $O(n)$ time a la Top. Ord.*
- 6: Obtain G_{rev} . ▷ *This takes $O(n + m)$ time.*
- 7: Run DFS(G_{rev}, π) and return the connected components of the forest F .

Theorem 2. The STRONGCONNCOMP algorithm returns the strongly connected components of G in a topological order of G^{scc} , in $O(n + m)$ time.

Proof. The running time is easiest to figure out. There are two DFSs which take $O(n + m)$ time. One needs to sort the last in decreasing order. As in TOPOLOGICAL ORDERING, this can be done in $O(n)$ time. What is interesting is the correctness of the algorithm.

Let C_1, \dots, C_k be the components of G^{scc} . We claim that the components of the final forest F returned in [Line 7](#) are these components returned in topological order of G^{scc} . More precisely, at the end of [Line 7](#), the variable fcomp contains the number of strongly connected components, that is, k , and for every $1 \leq t \leq \text{fcomp}$, the vertices $C_t := \{v : \text{Fcomp}[v] = t\}$ forms the t th strongly connected component, and the ordering C_1, \dots, C_k is the topological order of G^{scc} .

The discussion before the description of the algorithm indicates that C_1 is a sink component of $(G_{\text{rev}})^{\text{scc}}$ (to remind, the first vertex in π must lie in the *source* component of G^{scc} which is the *sink* component of $(G_{\text{rev}})^{\text{scc}}$). That is, C_1 is a true strongly connected component, and being the sink component of $(G_{\text{rev}})^{\text{scc}}$ it is indeed a source component of G^{scc} . What about the remaining C_i ’s?

To argue about them, we assume the statement is true for the first i components returned in [Line 7](#). That is, there are indeed the first i components in some topological order of G^{scc} . Next, we argue about the $(i + 1)$ th component.

Note, by the way DFS works, that the vertex v picked as root at the $(i + 1)$ th step of [Line 7](#) is the first vertex in π not in $C_1 \cup \dots \cup C_i$. Since π is the decreasing order of lasts obtained in [Line 4](#), we get that $\text{last}[v] = \max_{u \notin C_1, \dots, C_i} \text{last}[u]$. Suppose v lies in component C_j of G^{scc} . We now claim that [Lemma 2](#) implies C_j is a *source* component in $G^{\text{scc}} \setminus (C_1 \cup \dots \cup C_i)$. Suppose not: then there is some (C_k, C_j) edge in G^{scc} with $k > i$. From the lemma, we get that $f(C_k) > f(C_j)$ (otherwise the Lemma is violated). But this contradicts the choice of v . Therefore C_j is a source component of $G^{\text{scc}} \setminus (C_1 \cup \dots \cup C_i)$. That is, C_j is a sink component of $(G_{\text{rev}})^{\text{scc}} \setminus (C_1 \cup \dots \cup C_i)$. Now notice that the DFS run from vertex v on G_{rev} will only discover vertices in C_j as it is a sink component in $(G_{\text{rev}})^{\text{scc}} \setminus (C_1 \cup \dots \cup C_i)$. That is, the $(i + 1)$ th step discovers a component C_{i+1} which is a sink component of $(G_{\text{rev}})^{\text{scc}} \setminus (C_1 \cup \dots \cup C_i)$. That is, it is a source component of $G^{\text{scc}} \setminus (C_1 \cup \dots \cup C_i)$. But that implies C_{i+1} is the next component in the topological order of G^{scc} after C_1, \dots, C_i . \square