

CS 341: ALGORITHMS

Lecture 12: graph algorithms III – DAG testing, topsort, SCC
Readings: see website

Trevor Brown
<https://student.cs.uwaterloo.ca/~cs341>
trevor.brown@uwaterloo.ca

1

DFS APPLICATION: TESTING WHETHER A GRAPH IS A DAG

A directed graph G is a **directed acyclic graph**, or **DAG**, if G contains no directed cycle.

2

Lemma 6.7
A directed graph is a DAG if and only if a depth-first search encounters no back edges.

Proof.
(\Rightarrow): Any back edge creates a directed cycle.

Back edge:
points to an ancestor
in the DFS forest

3

- Case (\Leftarrow): Suppose \exists directed cycle. Show \exists back edge.
 - Let $v_1, v_2, \dots, v_k, v_1$ be a directed cycle
 - WLOG let v_1 be *earliest* discovered node in the cycle

edge type	discovery/finish times
tree	$d[v_k] < d[v_1] < f[v_1] < f[v_k]$
forward	$d[v_1] < d[v_2] < f[v_2] < f[v_1]$
back	$d[v_1] < d[v_k] < f[v_1] < f[v_k]$
cross	$d[v_1] < f[v_1] < d[v_2] < f[v_2]$

Discovered before v_2, \dots, v_k

Consider edge (v_k, v_1)

Since $d[v_1] < d[v_k]$, (v_k, v_1) must be a back or cross edge. Why?

So when v_1 is discovered, v_2, \dots, v_k are all white

Recall: every node v_i that is white-reachable from v_1 when we discover v_1 (call $DFSVisit(v_1)$) turns black before v_1 ($f[v_1] < f[v_i]$)

So v_k must turn black before v_1 , and we have $f[v_k] < f[v_1]$.

Thus, (v_k, v_1) must be a back edge. QED

4

TURNING THE LEMMA INTO AN ALGORITHM

Lemma 6.7
A directed graph is a DAG if and only if a depth-first search encounters no back edges.

- Search for back edges
- How to identify a back-edge?

edge type	colour of v	discovery/finish times
tree	white	$d[u] < d[v] < f[v] < f[u]$
forward	black	$d[u] < d[v] < f[v] < f[u]$
back	gray	$d[v] < d[u] < f[v] < f[u]$
cross	black	$d[v] < f[v] < d[u] < f[u]$

When we observe an edge from u to v , check if v is gray

Backedge

5

DFS: TESTING WHETHER A GRAPH IS A DAG

```

1 global variables:
2 pred[1..n] = [null, null, ..., null]
3 colour[1..n] = [white, white, ..., white]
4 d[1..n] = [0, 0, ..., 0] // discovery times
5 f[1..n] = [0, 0, ..., 0] // finish times
6 time = 0
7 DAG = true
8
9 IsDAG(adj[1..n])
10 for v = 1..n
11   if colour[v] == white
12     DFSVisit(adj, v)
13   return DAG
14
15 DFSVisit(adj[1..n], v)
16   colour[v] = gray
17   time = time + 1
18   d[v] = time
19
20   for each w in adj[v]
21     if colour[w] == white
22       pred[w] = v
23       DFSVisit(w)
24     if colour[w] == gray
25       DAG = false
26
27   colour[v] = black
28   time = time + 1
29   f[v] = time
    
```

6

EXAMPLE

Back edge found! So we set DAG = false

7

TOPOLOGICAL SORT

Finding node orderings that satisfy given constraints

8

DEPENDENCY GRAPH

- Edge $\{u, v\}$ means u must be completed **before** v

Example problem: getting dressed in the morning

Watch any time

Pants before belt

Socks before shoes

Could do various things first. Which ones are possible? What do they have in common?

9

Try to order nodes linearly so there are only pointers from left to right!

Possible IFF graph is a DAG

Topological sort

10

FORMAL DEFINITION

A directed graph $G = (V, E)$ has a **topological ordering**, or **topological sort**, if there is a linear ordering $<$ of all the vertices in V such that $u < v$ whenever $uv \in E$.

Graph G

Topological sort of G

$v_2 < v_5 < v_4 < v_1 < v_6$

Edges are directed only left-to-right in this ordering

11

USEFUL FACT

Lemma 6.5
A DAG contains a vertex of indegree 0.

Proof.
Suppose we have a directed graph in which every vertex has positive indegree. Let v_1 be any vertex. For every $i \geq 1$, let $v_{i+1}v_i$ be an arc. In the sequence v_1, v_2, v_3, \dots , consider the first repeated vertex, $v_i = v_j$ where $j > i$. Then $v_i, v_{j-1}, \dots, v_i, v_j$ is a directed cycle.

One of these must be repeated. So there is a cycle!

12

TOPOLOGICAL SORT VIA DFS

- We can implement topological sort by using **DFS!**
- The **finishing times** of nodes help us
- Understanding this algo will be **key** for understanding **strongly connected components**

Lemma 6.8

Suppose D is a DAG. Then $f[v] < f[u]$ for every arc uv .

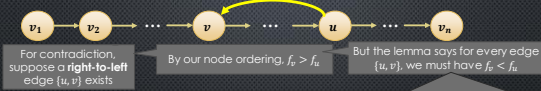
Recall from DAG-testing: there are **no back edges** in a DAG

edge type	colour of v	discovery/finish times
tree	white	$d[u] < d[v] < f[v] < f[u]$
forward	black	$d[u] < d[v] < f[v] < f[u]$
back	gray	$d[v] < d[u] < f[u] < f[v]$
cross	black	$d[v] < f[v] < d[u] < f[u]$



Theorem: if D is a DAG, and we order vertices in **reverse order of finishing time**, (i.e., by largest to smallest finish time) then we get a topological ordering!

To see **why**, suppose D is a DAG and we order nodes in this way, so $f_{v_1} > f_{v_2} > \dots > f_{v_{n-1}} > f_{v_n}$



For contradiction, suppose a **right-to-left edge** (u, v) exists

By our node ordering, $f_v > f_u$

But the lemma says for every edge (u, v) , we must have $f_u < f_v$

Lemma 6.8
Suppose D is a DAG. Then $f[v] < f[u]$ for every arc uv .

Contradiction! Right-to-left edge cannot exist. So is is a topological ordering.

TOPOLOGICAL ORDERING VIA DFS $O(n + m)$ w/adj. lists

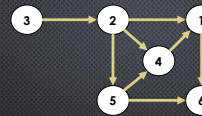
```

1 global variables:
2 pred[1..n] = [null, null, ..., null]
3 colour[1..n] = [white, white, ..., white]
4 d[1..n] = [0, 0, ..., 0] // discovery times
5 f[1..n] = [0, 0, ..., 0] // finish times
6 time = 0
7 DAG = true
8
9 TopologicalSort(adj[1..n])
10 S = new stack
11 for v = 1..n
12   if colour[v] == white
13     DFSVisit(adj, v, S)
14   if DAG then return S
15   return null
16
17 DFSVisit(adj[1..n], v, S)
18 colour[v] = gray
19 time = time + 1
20 d[v] = time
21
22 for each w in adj[v]
23   if colour[w] == white
24     pred[w] = v
25     DFSVisit(w)
26   if colour[w] == gray
27     DAG = false
28   colour[v] = black
29
30 S.push(v)
31 time = time + 1
32 f[v] = time
    
```

Push smallest finishing time first → pop largest first

Save each node when it finishes

HOME EXERCISE: RUN ON THIS GRAPH

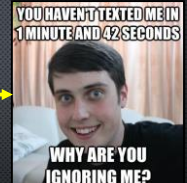


The initial calls are $DFSvisit(1)$, $DFSvisit(2)$ and $DFSvisit(3)$.

The discovery/finish times are as follows:

v	$d[v]$	$f[v]$	v	$d[v]$	$f[v]$
1	1	4	4	6	7
2	5	10	5	8	9
3	11	12	6	2	3

The topological ordering is 3, 2, 5, 4, 1, 6 (reverse order of finishing time).



STRONGLY CONNECTED COMPONENTS

STRONGLY CONNECTED COMPONENTS

- This graph could be divided into **two graphs** that are each strongly connected

These are called **strongly connected components (SCCs)**

19

STRONGLY CONNECTED COMPONENTS

- It could also be divided into **three graphs**... Not maximal

- But we want our SCCs to be **maximal** (as large as possible)

20

STRONGLY CONNECTED COMPONENTS

- So, the goal is to find **these** (maximal) SCCs:

21

APPLICATIONS OF SCCs AND COMPONENT GRAPHS

- Finding **all cyclic** dependencies in code
- Can find **single cycle** with an easier DFS-based algorithm
- But it is nicer to find **all cycles** at once, so you don't have to fix one to expose another

22

APPLICATIONS OF SCCs AND COMPONENT GRAPHS

- Data filtering** before running other algorithms
- maps; nodes = intersections, edges = roads
- Don't want to run path finding algorithm on the entire **global** graph!
- Throw away everything except the (maximal) SCC containing source & target

23

COMPONENT GRAPH

Consider this graph

These are its SCCs

Can there be a **cycle** in the component graph?

No! If there are paths both ways between components, they are actually **the same SCC**

Component graph is a **DAG!**

It has one node for each SCC

And an edge between two nodes IFF there is an edge between the corresponding SCCs

24

BRAINSTORMING AN ALGORITHM

- What if we run DFS, then reverse all edges, then run DFS (like checking whether an entire graph is strongly connected?)

This will definitely visit every node in a 's SCC
And in fact it might visit other SCCs as well...

Showing discovery times

Showing finish times

25

- What if we run DFS, then reverse all edges, then run DFS?

We fail to identify SCC (h, i)

Problem: from h, we can reach other SCCs

What if we perform DFSVisit calls in a different order in the reverse graph?

Other reachable SCCs should be visited first

Then, each DFSVisit will visit exactly one SCC (So we don't visit them again)

26

Consider component graph C_G of G (which we want to compute)

If we call DFSVisit in G from largest to smallest finish times, we can reach other SCCs.

graph G

graph H (reversed edges)

Recall lemma: edge uv in DAG implies $f(u) > f(v)$

27

SCC ALGORITHM

```

1 SCC(adj[1..n])
2 DFS(adj)
3   let order[1..n] = node labels sorted by
4     largest to smallest finish time
5
6   reverse all edges in adj
7
8   colour[1..n] = [white, ..., white]
9   comp[1..n] = [0, ..., 0]
10  for i = 1..n
11    v = order[i]
12    if colour[v] == white
13      scc = scc + 1
14      SCCVisit(adj, v, scc, colour, comp)
15
16  return comp
    
```

18 SCCVisit(adj[1..n], v, scc, colour, comp)

19 colour[v] = gray

20 comp[v] = scc

21 for each w in adj[v]

22 if colour[w] == white

23 SCCVisit(w)

24

25 colour[v] = black

26

This is called Sharir's algorithm (sometimes Kosaraju's algorithm). This paper first introduced it.

28

Running Sharir's Algorithm

Phase 1: DFS to get finish times

Phase 2: DFSVisit reverse graph by reverse finish times

scc is shown

scc = 4

29

TIME COMPLEXITY?

```

1 SCC(adj[1..n])
2 DFS(adj)
3   let order[1..n] = node labels sorted by
4     largest to smallest finish time
5
6   reverse all edges in adj
7
8   colour[1..n] = [white, ..., white]
9   comp[1..n] = [0, ..., 0]
10  for i = 1..n
11    v = order[i]
12    if colour[v] == white
13      scc = scc + 1
14      SCCVisit(adj, v, scc, colour, comp)
15
16  return comp
    
```

$O(n+m)$

$O(n)$

$O(n+m)$

$O(n)$

18 SCCVisit(adj[1..n], v, scc, colour, comp)

19 colour[v] = gray

20 comp[v] = scc

21 for each w in adj[v]

22 if colour[w] == white

23 SCCVisit(w)

24

25 colour[v] = black

26

Total of $O(n+m)$ work over all n iterations of the i loop

(each edge is inspected once, each node is visited once, constant work per visited node/inspected edge)

Total $O(n+m)$

30

CORRECTNESS

- Want to prove that each top-level call to SCCVisit explores **exactly** the nodes in **one SCC**
- Proof hinges on a key lemma that talks about the **finish times of SCCs** in the **component graph**
- To talk about finish times of **SCCs**, we need a definition...

31

A KEY DEFINITION

For a strongly connected component C , let $d[C] = \min\{d[v] : v \in C\}$ and $f[C] = \max\{f[v] : v \in C\}$.

Examples of calculations:
 $d[C_1] = 19 = \min\{d[i], d[j], d[k]\}$
 $f[C_1] = 24 = \max\{f[i], f[j], f[k]\}$

32

A KEY LEMMA

- Lemma:** if C_i, C_j are SCCs and there is an edge $C_i \rightarrow C_j$ in G , then $f[C_i] > f[C_j]$

- Proof.** Case 1 ($d[C_i] < d[C_j]$):

- Let u be the earliest discovered node in C_i
- All nodes in $C_i \cup C_j$ are white-reachable from u , so they are **descendants in the DFS forest** and **finish before u**
- So $f[C_i] = f[u] > f[C_j]$

33

A KEY LEMMA

- Lemma:** if C_i, C_j are SCCs and there is an edge $C_i \rightarrow C_j$ in G , then $f[C_i] > f[C_j]$

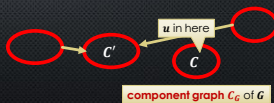
- Proof.** Case 2 ($d[C_j] < d[C_i]$):

- Since component graph is a DAG, there is **no path $C_j \rightarrow C_i$**
- Thus, **no nodes** in C_i are reachable from C_j
- So we discover C_j and finish C_j **without** discovering C_i
- Therefore $d[C_j] < f[C_j] < d[C_i] < f[C_i]$. **QED**

34

COMPLETING THE PROOF

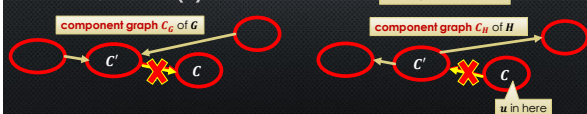
- Suppose we have performed DFS to get our finish times, and we are about to perform SCCVisits on the reverse graph
- We prove each top-level SCCVisit call visits precisely one SCC**
- Consider the first top-level SCCVisit(u)
- Let C be the SCC containing u and C' be any other SCC
- Since we call SCCVisit on nodes starting from the **largest finish time**,
 - We know $f(C) > f(C')$



35

COMPLETING THE PROOF

- We know $f(C) > f(C')$
- By Lemma: if there were an edge $C' \rightarrow C$ in G , then we would have $f(C') > f(C)$
 - So there is no edge $C' \rightarrow C$ in G
 - and hence **no edge $C \rightarrow C'$ in H**
 - So, **SCCVisit(u) in H cannot visit C'**



36

IF WE HAVE TIME

topological sort without relying on DFS

37

EXISTENCE OF A TOPOLOGICAL SORT ORDER

Theorem 6.6
A directed graph D has a topological sort if and only if it is a DAG.

Proof.
(\Rightarrow): Suppose D has a directed cycle $v_1, v_2, \dots, v_j, v_1$. Then $v_1 < v_2 < \dots < v_j < v_1$, so a topological ordering does not exist.
(\Leftarrow): Suppose D is a DAG. Then the algorithm below constructs a topological ordering.

38

```

1 Kahn(adj[1..n])
2   indeg[1..n] = [0, ..., 0]
3   for each edge (u,v) in adj
4     indeg[v] = indeg[v] + 1
5
6   order = new list
7   q = new queue containing {v : indeg[v] == 0}
8   for i = 1..n
9     if q.empty() return null
10    v = q.dequeue()
11    order.append(v)
12
13    for each w in adj[v]
14      indeg[w] = indeg[w] - 1
15      if indeg[w] == 0 then q.enqueue(w)
16
17   return order
    
```

$\text{indeg}[v]$ = # of edges pointing into node v = number of unsatisfied constraints on v

Nodes with $\text{indeg } 0$ have no unsatisfied dependencies

So this step is enqueueing nodes whose dependencies are already satisfied

q always contains nodes with no unsatisfied dependencies ($\text{indeg } 0$)

No such order!

Add v to the topological order

Remove v 's out edges. If we have now satisfied all dependencies for some w , add w to the queue also.

39

EXAMPLE (KAHN'S ALGORITHM)

Compute **indegree** for all vertices

For each node v
For each w in $\text{adj}[v]$
 $w.\text{deg} = w.\text{deg} + 1$

vertices with $\text{indeg } 0$ go into the queue

Until Q is empty: pop, output that element, decrement its neighbours, enqueue new $\text{indeg } 0$'s

Queue Q
3 2 5 4 1 6

Output

40

```

1 Kahn(adj[1..n])
2   indeg[1..n] = [0, ..., 0]
3   for each edge (u,v) in adj
4     indeg[v] = indeg[v] + 1
5
6   order = new list
7   q = new queue containing {v : indeg[v] == 0}
8   for i = 1..n
9     if q.empty() return null
10    v = q.dequeue()
11    order.append(v)
12
13    for each w in adj[v]
14      indeg[w] = indeg[w] - 1
15      if indeg[w] == 0 then q.enqueue(w)
16
17   return order
    
```

Running time with adjacency lists?

$O(n)$

$O(n+m)$ total work over all iterations

$O(n)$ iterations

$O(1)$ per check

$O(1)$

$\sum_{v \in V} \text{deg}(v) \in O(n+m)$

total work over all nodes v

Total $O(n+m)$

41

BONUS SLIDES

42

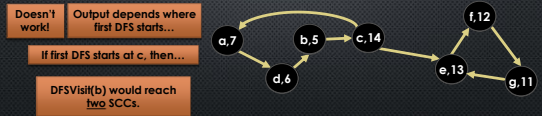
SCC: HOW ABOUT A DIFFERENT ORDERING?

- Rather than doing DFS in the **reverse** graph in order of **decreasing** finish times
- Why not do DFS in the **original** graph in order of **increasing** finish times?
- Exercise: does this work?

43

SCC: HOW ABOUT A DIFFERENT ORDERING?

- Why not do DFS in the **original** graph in order of **increasing** finish times?



44