

# Tutorial 1 - Graph Representation, Traversal and MST

COMP9418 – Advanced Topics in Statistical Machine Learning

Lecturer: Gustavo Batista

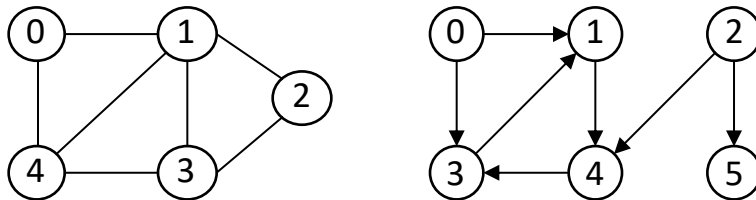
---

**Last revision:** Tuesday 8<sup>th</sup> September, 2020 at 15:25

In this tutorial, we will review some basic concepts of graph representation, traversal and minimum spanning trees (MST). In particular, we will look at some relevant problems we can solve with graph traversal and efficient implementations of MST algorithms. This material will not be covered in the course, but it is essential to understand and implement several algorithms of Probabilistic Graphical Models.

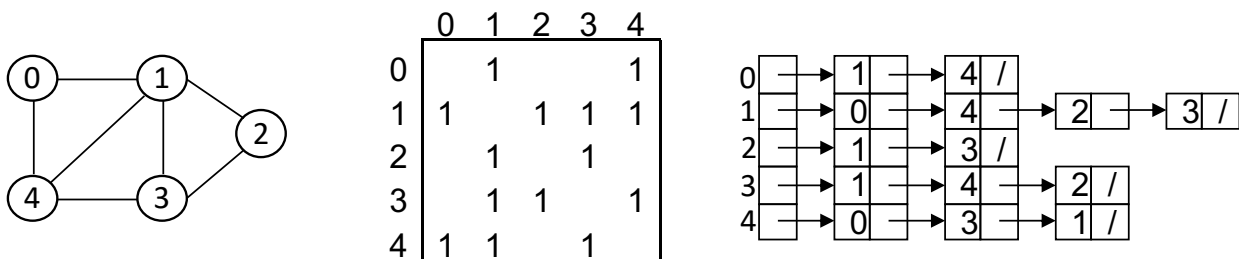
## Question 1 - Graph representation

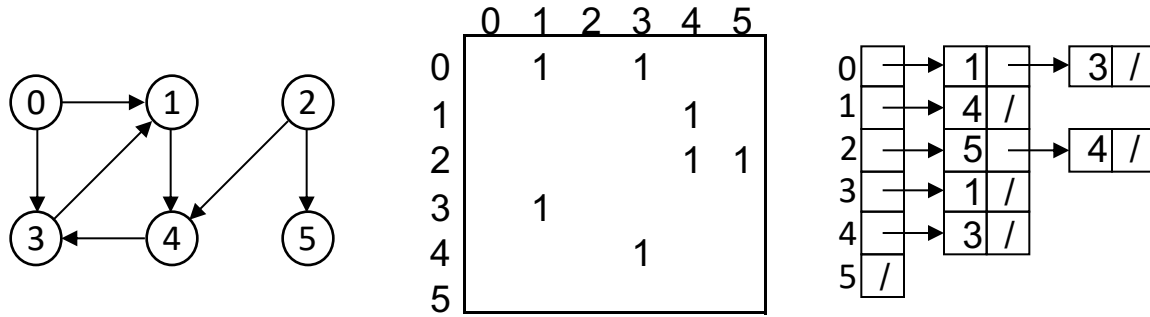
There are two frequently used representations for graphs: adjacency matrix and list. Both representations can be used for directed and undirected graphs, like the ones in the next figure.



Represent both graphs using the matrix and list representations. You do not need to write any code here, use a graphical representation such as boxes linked by arrows. Which of these two representations is typically preferred to implement operations over graphs? Why?

**Answer**





We usually prefer the adjacency list representation since it allows more efficient implementations of the basic operations on sparse graphs. For instance, for a graph with  $n$  vertices and  $m$  edges, graph traversal operations are implemented in  $O(n^2)$  operations with the matrix representation and  $O(n + m)$ . For sparse graph  $m \ll n^2$ .

## Question 2 - Graph traversal

Graph traversal is the process of visiting each vertex in a graph. There are two main variations, depth-first (DFS) and breadth-first search (BFS). Write the algorithms for each of these operations. Use a simplifying notation that allows you to write a single algorithm for both graph representations. For problems that can be solved by both algorithms, is there one search strategy (depth or breadth) that is preferred? Why?

**Answer**

DFS algorithm.

```

1 procedure DFS(G, v):
2   label v as grey
3   for all edges e in G.incidentEdges(v) do
4     w := G.adjacentVertex(v, e)
5     if vertex w is white then
6       DFS(G, w)
7   label v as black

```

BFS algorithm.

```

1 procedure BFS(G, v):
2   create a queue Q
3   enqueue v onto Q
4   mark v
5   while Q is not empty do
6     w := Q.dequeue()
7     for all edges e in G.adjacentEdges(w) do
8       x := G.adjacentVertex(w, e)
9       if x is not marked then
10        mark x
11        enqueue x onto Q

```

Both algorithms need a data structure to keep track of the search. DFS uses a stack and BFS a queue. Most implementations of DFS are recursive, so the algorithm avails of the program stack memory, making this data structure implicit. For problems that can be solved by both algorithms, we usually prefer DFS over BFS. The reason is the amount of memory necessary to run these search procedures, i.e., the maximum number of elements we can store in the stack or queue during the execution. For DFS, the memory complexity of  $O(h)$

where  $h$  is the maximum depth of the search tree (maximum depth of recursive calls). BFS has a memory complexity of  $O(n)$ . For many practical problems,  $h \ll n$ .

### Question 3 - Cycle detection in directed graphs

In the first part of this course, we will use directed acyclic graphs (DAGs) to represent statistical independencies between random variables. Therefore, it may be useful to develop an algorithm to detect cycles to verify if a directed graph is a DAG. The DFS algorithm can be adapted to detect the existence of cycles in a directed graph. A cycle is a non-empty path that links a vertex to itself. A well-known trick is to colour the vertices as we visit them. Vertices are all initially white and become grey as we first visit them. A grey vertex becomes black when we have visited all descendent nodes of this grey vertex, and we are ready to backtrack the recursion. Adapt your DFS algorithm to detect cycles in directed graphs using this colouring scheme. Is a cycle detected when we reach a white, grey or black node?

#### Answer

We will use the colouring vertex scheme of the DFS implementation to detect cycles. In this case, we stop when we identify the first cycle as an edge that links to a grey node. In directed graphs, such an arc is called a “back” edge. Notice that edges that connect to black nodes are “cross” edges and are not indicative of cycles.

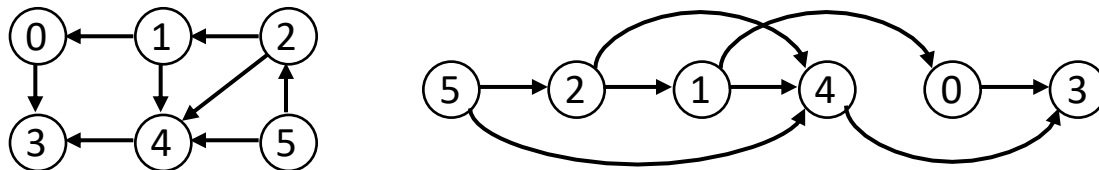
```

1 procedure findCycle(G, v):
2   label v as grey
3   for all edges e in G.incidentEdges(v) do
4     w := G.adjacentVertex(v, e)
5     if vertex w is white then
6       if findCycle(G, w) then
7         return true
8     else if vertex w is grey then
9       return true
10  label v as black
11  return false

```

### Question 4 - Topological sorting on DAGs

Many algorithms in this course require iterating through nodes, where all parents need to be processed before their children. For these algorithms, we use an ordering known as *topological ordering*. Topological sort or topological ordering of a DAG is a linear ordering of the vertices. For every directed edge from vertex  $u$  to vertex  $v$ ,  $u$  comes before  $v$  in the order. The next figure shows a topological ordering for a DAG.



Adapt the DFS algorithm to produce a topological order of a DAG. Do you need any additional data structure to report the results in proper order?

#### Answer

The idea is that a black node has no further outgoing edges to be explored. Consequently, all nodes that “depend” on this black vertex are already in the output stack. We use a stack to revert the order of the nodes.

```

1 procedure topologicalSort(G, v):
2   label v as grey
3   for all edges e in G.incidentEdges(v) do
4     w := G.adjacentVertex(v, e)
5     if vertex w is white then
6       topologicalSort(G, w)
7   label v as black
8   push v in s           # s is a stack with topological order

```

## Question 5 - Strongly connected directed graphs

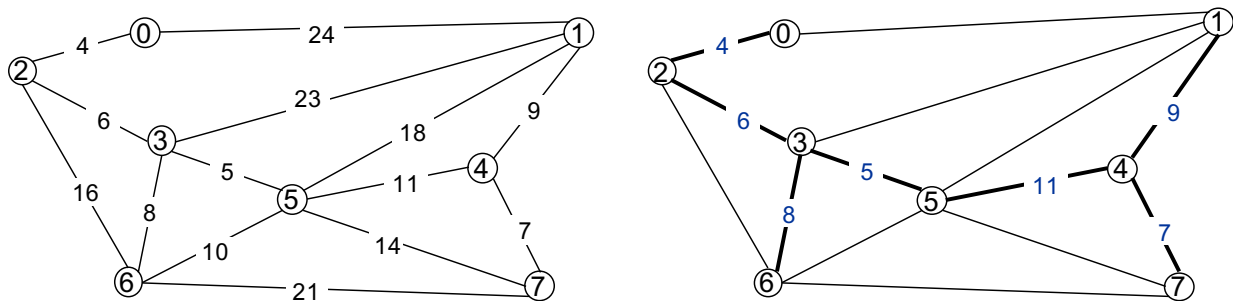
A directed graph is called strongly connected if there is a path in each direction between each pair of vertices. We can easily (but inefficiently) test if a directed graph is strongly connected by running  $n$  DFS searches from each vertex and confirming that we can reach every node in each search. However, according to Kosaraju's algorithm, two searches suffice, the first one in the graph  $G$  and the second one in its transpose  $G^T$ . Explain the intuition behind this algorithm.

### Answer

First,  $G^T$  is obtained by inverting the direction of each edge of  $G$ . For the matrix representation, we can get  $G^T$  using the matrix transpose operation. A DFS search from a node  $v$  that reaches all vertices in  $G$  confirms the existence of a directed path between  $v$  and every node of  $G$ . A DFS search from the same node  $v$  that reaches all vertices in  $G^T$  indicates that there is a directed path between all these nodes and  $v$ . Therefore,  $G$  is strongly connected since every node can reach every other node with a directed path that passes through  $v$ .

## Question 6 - Minimal spanning trees

A common operation over undirected graphs is known as a minimum spanning tree (MST). Given an undirected graph  $G$ , MST consists of finding a subset of the edges that forms a tree and includes every vertex of  $G$ . The sum of the edges' weights in the tree should be minimal. The next figure shows an example of a graph (left) and its minimum spanning tree (right).



Notice that for a graph with  $n$  vertices, the spanning tree must have exactly  $n - 1$  edges. Adding one more edge would create a cycle, and removing one edge would disconnect the tree creating a forest.

There are several algorithms to compute an MST of a graph, including the well-known Prim and Kruskal algorithms. These algorithms can be trivially adapted to calculate the *maximum* spanning tree by selecting the edges with the highest cost instead of the smallest cost.

Write the pseudo-code of an MST algorithm such as PRIM. What is the time complexity of your algorithm? Which data structures are necessary to achieve such complexity?

## Answer

The PRIM pseudo-code is the following.

```
1 procedure Prim(G, s):
2   S := {s}
3   Q := empty priority queue
4   for all edges e in G.incidentEdges(s) do
5     Q.insert([e.cost, e.from, e.to])           # e.from = s
6
7   while (not empty(Q))
8     [cost, v, u] := Q.remove()                 # edge with smallest cost in Q
9     if u not in S
10      S := S U {u}
11      for all edges e in G.incidentEdges(u) do
12        if e.to not in S
13          Q.insert([e.cost, e.from, e.to])     #e.from = u
```

This algorithm has a  $O(m \log m)$  time complexity, where  $m$  is the number of edges. To achieve such time complexity, we had to use a priority queue data structure that allows to insert and remove elements with  $O(\log n)$  cost.