# Practical 4: NCBI Taxonomy

Ximena Moure  Eliya Tiram

10/10/2023, submission deadline 16/10/2023

1. **Given a file nodes.dmp for the NCBI taxonomy, write a Python script to store it in a rooted tree. Give the code of your Python script as your answer to this question**

   The rooted tree we have stored the data in is in the following python code. Since we build a rooted tree, as requested in the question, we ignore self loops, for example 1—1 in the data.

```python
import networkx as nx

# Define a function to process the data and construct the graph
def construct_tree_from_file(file_path):
    # Create a directed graph
    graph = nx.DiGraph()

    with open(file_path, 'r') as file:
        for line in file:
            columns = line.strip().split('|')
            if len(columns) >= 2:
                parent, child = map(int, columns[:2])

                # Skip self-loops (edges where parent and child are the same)
                if parent != child:
                    # Add nodes to the graph
                    graph.add_node(parent)
                    graph.add_node(child)

                    # Add an edge to represent the parent-child relationship
                    graph.add_edge(parent, child)

    return graph

# Path to the input file
file_path = "nodes.dmp"

# Construct the tree from the data in the file
tree = construct_tree_from_file(file_path)
```

Listing 1: Python code ex1

2. **Is the NCBI taxonomy a rooted tree, a directed acyclic graph, or a directed graph with cycles?**

   The NCBI taxonomy is a directed graph with cycles. In this section we do not ignore self loops so we will know the NCBI taxonomy type.

```python
import networkx as nx

# Define a function to process the data and construct the graph
def construct_tree_from_file(file_path):
    # Create a directed graph
    graph = nx.DiGraph()

    with open(file_path, 'r') as file:
        for line in file:
            columns = line.strip().split('|')
            if len(columns) >= 2:
                parent, child = map(int, columns[:2])
                # Add nodes to the graph
                graph.add_node(parent)
                graph.add_node(child)
                # Add an edge to represent the parent-child relationship
                graph.add_edge(parent, child)
```

```
18
19      return graph
20
21  # Path to the input file
22  file_path = "nodes.dmp"
23
24  # Construct the tree from the data in the file
25  tree = construct_tree_from_file(file_path)
26
27  # Check if it's a rooted tree, DAG, or graph with cycles
28  if nx.is_tree(tree):
29      print("The NCBI taxonomy is a rooted tree.")
30  elif nx.is_directed_acyclic_graph(tree):
31      print("The NCBI taxonomy is a directed acyclic graph (DAG).")
32  else:
33      print("The NCBI taxonomy is a directed graph with cycles.")
34
```

Listing 2: Python code ex2

3. **How many nodes are there in the NCBI taxonomy?**
   In this section we ignore self loops again and we get the number of nodes in the NCBI taxonomy is 2442791 and the number of edges in the NCBI taxonomy is 2442790.

```
1  # Number of nodes and edges in the NCBI taxonomy
2  num_nodes = len(tree.nodes())
3  num_edges = tree.number_of_edges()
4  print(f"The number of nodes in the NCBI taxonomy is {num_nodes}.")
5  print(f"The number of edges in the NCBI taxonomy is {num_edges}.")
6
```

Listing 3: Python code ex3

4. **Write a Python script to restrict the rooted tree to the seven standard taxonomic ranks: kingdom, phylum, class, order, family, genus, species.**

```
1  import networkx as nx
2
3  # Define the valid ranks
4  valid_ranks
5      = ["kingdom", "phylum", "class", "order", "family", "genus", "species"]
6
7  def construct_tree_from_file(file_path):
8      # Create a directed graph
9      graph = nx.DiGraph()
10
11     # Read the file and build the entire graph
12     with open(file_path, 'r') as file:
13         for line in file:
14             columns = line.strip().split('\t|\t')
15             if len(columns) >= 3:
16                 node_id, parent_id, rank = map(str.strip, columns[:3])
17                 node_id, parent_id = map(int, [node_id, parent_id])
18                 # Add nodes to the graph with rank attribute
19                 graph.add_node(node_id, rank=rank)
20                 graph.add_node(parent_id)
21                 # Add an edge to represent the parent-child relationship
22                 graph.add_edge(parent_id, node_id)
23
24     # Identify nodes to be removed and nodes to be retained
25     not_valid_nodes = []
26     for node, attrs in graph.nodes(data=True):
27         # Check if the rank of the current node is not in the valid ranks
28         if attrs['rank'] not in valid_ranks:
29             # If the rank is not valid, add the node to the not_nodes list
```

```
30              not_valid_nodes.append(node)
31
32       restricted_nodes = [node
         for node, attrs in graph.nodes(data=True) if attrs['rank'] in valid_ranks]
33
34       # Create a copy of the original graph for modification
35       rest_tax = graph.copy()
36
37       # Lists of successors and predecessors for the nodes which will be removed
38       successors = [list(graph.successors(i)) for i in not_valid_nodes]
39       predecessors = [list(graph.predecessors(i)) for i in not_valid_nodes]
40
41       # Add edges from the predecessor to the successors of nodes to be removed
42       rest_tax.add_edges_from(
43           [(predecessors[i][0], successors[
     i][j]) for i in range(len(successors)) for j in range(len(successors[i]))])
44
45       # Restrict the graph to nodes of standard ranks
46       rest_tax = nx.subgraph(rest_tax, restricted_nodes)
47
48       return rest_tax  # Return the restricted graph and the root node
49
50
51 # Path to the input file
52 file_path = "nodes.dmp"
53 restricted_taxonomy = construct_tree_from_file(file_path)
54 num_nodes = restricted_taxonomy.number_of_nodes()
55 num_edges = restricted_taxonomy.number_of_edges()
56 print(f"The number of nodes in the restricted taxonomy is {num_nodes}.")
57 print(f"The number of edges in the restricted taxonomy is {num_edges}.")
58
```
Listing 4: Python code ex4

5. **What is the name of the kingdom taxonomic rank in the NCBI taxonomy?**
   Superkingdom

6. **How many nodes are there in the NCBI taxonomy, once restricted to the seven standard taxonomic ranks?**
   There are 2112007 nodes.