

Practical 5: LCA Skeleton Tree

Ximena Moure Eliya Tiram

17/10/2023, submission deadline 23/10/2023

The following code is improved from the last delivery in order to get the restricted taxonomy to be able to do this practical. Now we do get 2111999 nodes.

```
1 import networkx as nx
2
3
4 def load_graph_from_file(file_path):
5     G = nx.DiGraph() # Initialize directed graph
6     with open(file_path, 'r') as file:
7         for line in file:
8             child, parent, rank = map(str.strip, line.split('|')[:3])
9             child, parent = int(child), int(parent)
10
11             # Skip self-loops
12             if child != parent:
13                 # Add edge
14                 G.add_edge(parent, child)
15                 # Assign rank attribute to the child node
16                 G.nodes[child]['rank'] = rank
17     return G
18
19
20 file_path = "nodes.dmp"
21
22 initial_tree = load_graph_from_file(file_path)
23
24 valid_ranks
25     = ["superkingdom", "phylum", "class", "order", "family", "genus", "species"]
26 # Filter nodes based on rank attribute
27 nodes_to_remove = [node for node, attrs
28                     in initial_tree.nodes(data=True) if attrs.get('rank') not in valid_ranks]
29 # Create a copy of the graph
30 restricted_taxonomy = initial_tree.copy()
31
32 def create_restricted_tax(graph, nodes_to_delete):
33     for node in nodes_to_delete:
34         successors = list(graph.successors(node))
35         predecessors = list(graph.predecessors(node))
36         if node != 1:
37             graph.remove_node(node)
38             for source in predecessors:
39                 for target in successors:
40                     if source != target and not graph.has_edge(source, target):
41                         graph.add_edge(source, target)
42
43 create_restricted_tax(restricted_taxonomy, nodes_to_remove)
44 print('Tree: ', nx.is_tree(restricted_taxonomy))
45 print('
46     Rooted:', nx.is_arborescence(restricted_taxonomy)) # to check if it is rooted
47 print('Number
48     of nodes in the restricted taxonomy: ', restricted_taxonomy.number_of_nodes())
49 print('Number
50     of edges in the restricted version: ', restricted_taxonomy.number_of_edges())
```

Listing 1: Python code to extract restricted taxonomy

1. Given a file nodes.dmp for the NCBI taxonomy and a file mapping.txt of mappings of sequence reads to

NCBI taxonomic identifiers, write a Python script to find the lineage of the mapped nodes for each sequence read. Give the code of your Python script as your answer to this question, using the LATEX package listings.

```

1 def find_lineage(taxonomy_graph, node, cache=None):
2     if cache and node in cache:
3         return cache[node]
4
5     lineage = []
6     current = node
7     while current in taxonomy_graph:
8         lineage.append(current)
9         parents = list(taxonomy_graph.predecessors(current))
10        if not parents:
11            break
12        current = parents[0]
13
14    lineage.reverse()
15
16    if cache is not None:
17        cache[node] = lineage
18
19    return lineage
20
21
22 # Reading the mapping file
23 mapping = {}
24 with open('mapping.txt', 'r') as file:
25     for line in file:
26         li = line.strip().split()
27         sq_read_id = li[0]
28         nodes = [int(node) for node in li[1:]]
29         mapping[sq_read_id] = nodes
30
31 lineages = {}
32 lineage_cache = {}
33 for read_id, taxonomic_nodes in mapping.items():
34     lineage_for_read = [find_lineage
35                         (restricted_taxonomy, node, lineage_cache) for node in taxonomic_nodes]
36     lineages[read_id] = lineage_for_read
37
38 print("lineages", lineages)

```

Listing 2: Python code ex1

2. **What is the taxonomic rank of the LCA of the mapped nodes for each of the sequence reads?**

The taxonomic rank of the LCA of the mapped nodes for each of the sequence reads is the root.

3. **Do the sequence reads come from archaea, bacteria, eukaryota, or viruses?**

The answer is:

R00010: eukaryota: 8, bacteria: 2

R00020: eukaryota: 16, viruses: 1, bacteria: 3

R00030: eukaryota: 21, bacteria: 8, archaea: 1

R00040: eukaryota: 26, bacteria: 14

R00050: bacteria: 7, eukaryota: 42, viruses: 1

R00060: viruses: 3, eukaryota: 43, bacteria: 14

R00070: bacteria: 21, eukaryota: 49

R00080: others: 1, eukaryota: 57, bacteria: 20, viruses: 2

R00090: others: 1, eukaryota: 59, archaea: 2, bacteria: 21, viruses: 7

R00100: others: 1, eukaryota: 79, bacteria: 19, archaea: 1

Thus, not all come from archaea, bacteria, eukaryota, or viruses. Most of them come from eukaryota.

```
1 from collections import defaultdict
2
3 sk_matching
4   = { 2: 'bacteria', 10239: 'viruses', 2157: 'archaea', 2759: 'eukaryota'}
5 superkingdoms = {}
6
7 for read_id, seq in lineages.items():
8     sk_read = defaultdict(int)
9     sk_read['others'] = 0
10    for l in seq:
11        sk_read[sk_matching.get(l[1], 'others')] += 1
12
13    superkingdoms[read_id] = sk_read
14
15 for read_id, sk_counts in superkingdoms.items():
16     formatted_output = ', '.join([f"{k}: {v}" for k, v in dict(sk_counts).items() if not (k == 'others' and v == 0)])
17     print(f"{read_id}: {formatted_output}")
```

Listing 3: Python code ex3

4. Given a file nodes.dmp for the NCBI taxonomy and a file mapping.txt of mappings of sequence reads to NCBI taxonomic identifiers, write a Python script to build the LCA skeleton tree for each sequence read. Give the code of your Python script as your answer to this question, using the LATEX package listings.

The following code uses the lineages from exercise 1

```
1 skeleton_trees = {}
2
3 for read_id, seq_lineages in lineages.items():
4     all_nodes = {node for lineage in seq_lineages for node in lineage}
5     skeleton_graph = nx.DiGraph()
6     skeleton_graph.add_nodes_from(all_nodes)
7
8     for lineage in seq_lineages:
9         edges_from_lineage = zip(lineage[:-1], lineage[1:])
10        skeleton_graph.add_edges_from(edges_from_lineage)
11
12    skeleton_trees[read_id] = skeleton_graph
13
14    print(f"{read_id} has {skeleton_trees[read_id].number_of_nodes()} nodes")
15
```

Listing 4: Python code ex4

5. How many nodes are there in the LCA skeleton tree for each of the sequence reads?

The answer to this question is in the code provided for question number 4 in the following line (which is inside the loop, so it is printed for every read)

```
1     print(f"{read_id} has {skeleton_trees[read_id].number_of_nodes()} nodes")
2
```

Listing 5: Python code ex5

R00010 has 47 nodes
R00020 has 85 nodes
R00030 has 130 nodes
R00040 has 154 nodes
R00050 has 201 nodes
R00060 has 214 nodes

R00070 has 266 nodes
R00080 has 283 nodes
R00090 has 305 nodes
R00100 has 344 nodes

6. What is the taxonomic rank of the root of the LCA skeleton tree for each of the sequence reads?

The taxonomic rank of the root of the LCA skeleton tree for each of the sequence reads is the root