

## Practical 6: LCA Mapping

Ximena Moure Eliya Tiram

24/10/2023, submission deadline 30/10/2023

This code is used to build the restricted tree which will be used in the exercises of this lab.

```
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
18     return G
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 LCA mapping 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
32 lineages = {}
33 lineage_cache = {}
34 for read_id, taxonomic_nodes in mapping.items():
35     lineage_for_read = [find_lineage
36                         (restricted_taxonomy, node, lineage_cache) for node in taxonomic_nodes]
37     lineages[read_id] = lineage_for_read
38
39 # LCAs for each read
40 LCAs = {}
41
42 for read_id, seq in lineages.items():
43     min_length = min(len(lst) for lst in seq)
44     lst1 = seq[0]
45
46     i = 0
47     while i < min_length:
48         current_item = lst1[i]
49         # Check if the current item is the same for all sequences
50         all_same = all(lst[i] == current_item for lst in seq[1:])
51
52         # If any sequence differs at this position, break
53         if not all_same:
54             break
55
56         i += 1
57
58     # Store the last common ancestor (i.e., the item before the divergence)
59     LCAs[read_id] = lst1[i - 1] if i > 0 else None

```

```

61 print("The LCAs for each read: \n", LCAs)
62
63 # Retrieve the rank for each LCA
64 lca_ranks = {read_id: restricted_taxonomy.
65               nodes[lca].get('rank') for read_id, lca in LCAs.items() if lca is not None}
66
67 highest_rank = None
68 lowest_rank = None
69
70 # Determine the highest rank
71 for rank in lca_ranks.values():
72     if highest_rank
73         is None or valid_ranks.index(rank) < valid_ranks.index(highest_rank):
74             highest_rank = rank
75
76 print("Highest taxonomic rank for the LCAs:", highest_rank)
77
78 # Retrieve the rank for each LCA if it's not None
79 lca_ranks = {read_id: restricted_taxonomy.
80               nodes[lca].get('rank') for read_id, lca in LCAs.items() if lca is not None}
81
82 # Determine the lowest rank
83 for rank in lca_ranks.values():
84     if lowest_rank
85         is None or valid_ranks.index(rank) > valid_ranks.index(lowest_rank):
86             lowest_rank = rank
87
88 print("Lowest taxonomic rank for the LCAs:", lowest_rank)

```

Listing 2: Python code ex1

Result is shown in the image below (Figure: 1 ).

The LCAs for each read:

```
{'R00010': 1, 'R00020': 1, 'R00030': 1, 'R00040': 1, 'R00050': 1, 'R00060': 1, 'R00070': 1, 'R00080': 1, 'R00090': 1, 'R
```

Figure 1: Ex1 result

2. What is the highest taxonomic rank (that is, toward kingdom) for these LCA mappings?

It is no rank.

3. What is the lowest taxonomic rank (that is, toward species) for these LCA mappings?

It is no rank

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 find the optimal (in terms of the F-measure) taxonomic assignment for each sequence read. Give the code of your Python script as your answer to this question, using the LATEX package listings

In order to do this part, we modified the code to get the skeleton tree that we delivered for the previous lab. The new code is the following

```
1 def build_skeleton_tree(lineages):
2     skeleton_graph = nx.DiGraph()
3
4     for lineage in lineages:
5         # Add nodes and edges from the lineage to the skeleton graph
```

```

6         skeleton_graph.add_nodes_from(lineage)
7         skeleton_graph.add_edges_from(zip(lineage[:-1], lineage[1:]))
8
9     return skeleton_graph
10
11
12 skeleton_trees = {}
13
14 for read_id, seq_lineages in lineages.items():
15     skeleton_trees[read_id] = build_skeleton_tree(seq_lineages)
16
17
18 def compressLinearPaths(tree):
19     """
20     Contract elementary paths in the tree.
21     Elementary
22     paths are paths where a node has only one child and can be contracted.
23     This function traverses the tree in a bottom-up manner.
24     """
25     contracted_tree = tree.copy()
26     nodes_to_process = [node for node in nx.
27         topological_sort(contracting_tree) if contracting_tree.out_degree(node) == 1]
28
29     for current_node in nodes_to_process:
30         parent = next(contracting_tree.predecessors(current_node), None)
31         child = next(contracting_tree.successors(current_node), None)
32
33         # Skip if it's a root or leaf node
34         if parent is None or child is None:
35             continue
36
37         # Contract the node by connecting its parent to its child
38         contracting_tree.add_edge(parent, child)
39         contracting_tree.remove_node(current_node)
40
41     return contracting_tree
42
43
44 contracting_skel_trees = {}
45 # Counting nodes after contracting skeleton trees
46 nodes_per_read = {}
47 for read_id, lca_tree in skeleton_trees.items():
48     contracting_skel_tree = compressLinearPaths(lca_tree)
49     contracting_skel_trees[read_id] = contracting_skel_tree
50     num_nodes = len(contracting_skel_tree.nodes)
51     nodes_per_read[read_id] = num_nodes
52
53
54 print("Details of each skeleton tree:")
55 for read_id, tree in contracting_skel_trees.items():
56     num_nodes = tree.number_of_nodes()
57     num_edges = tree.number_of_edges()
58     print(f"Read
ID: {read_id}, Number of Nodes: {num_nodes}, Number of Edges: {num_edges}")

```

Listing 3: Python code for building the LCA skeleton tree for each sequence read

```

Details of each skeleton tree:
Read ID: R00010, Number of Nodes: 16, Number of Edges: 15
Read ID: R00020, Number of Nodes: 32, Number of Edges: 31
Read ID: R00030, Number of Nodes: 44, Number of Edges: 43
Read ID: R00040, Number of Nodes: 58, Number of Edges: 57
Read ID: R00050, Number of Nodes: 70, Number of Edges: 69
Read ID: R00060, Number of Nodes: 82, Number of Edges: 81
Read ID: R00070, Number of Nodes: 99, Number of Edges: 98
Read ID: R00080, Number of Nodes: 106, Number of Edges: 105
Read ID: R00090, Number of Nodes: 129, Number of Edges: 128
Read ID: R00100, Number of Nodes: 135, Number of Edges: 134

```

Figure 2: Skeleton tree for each sequence read

Now, we can do exercise 4. The code is the following

```

1
2 def calculate_precision_recall(tp, fp, fn):
3     precision = tp / (tp + fp) if (tp + fp) > 0 else 0
4     recall = tp / (tp + fn) if (tp + fn) > 0 else 0
5     return precision, recall
6
7
8 def calculate_f_measure(precision, recall):
9     return 2/(1/precision + 1/recall)
10
11
12 # Optimal calculation for each sequence read
13 results = []
14 optimal_ranks = {}
15
16 for read_id in mapping.keys():
17     # Retrieve the contracted skeleton tree for the read
18     tree = contracted_skel_trees.get(read_id, None)
19
20     # Check if the tree exists
21     if tree is None:
22         print(f"No skeleton tree found for read {read_id}")
23         continue
24
25     max_node = ''
26     max_F = 0
27     max_rank = ''
28     for node in tree.nodes():
29         # Calculate F_measure for each internal node
30         if node not in mapping[read_id]:
31             descendants_of_node = set(nx.descendants(tree, node))
32
33             true_positives = len(descendants_of_node.intersection(mapping[read_id]))
34             false_negatives = len(mapping[read_id]) - true_positives
35             descendants = 0
36             for d in set(nx.descendants(restricted_taxonomy, node)):
37                 if restricted_taxonomy.out_degree(d) == 0:
38                     descendants += 1
39             false_positives = descendants - true_positives
40             precision, recall = calculate_precision_recall(
41                 true_positives, false_positives, false_negatives)
42             f_measure = calculate_f_measure(precision, recall)
43
44     # Get the max F_measure

```

```

43         if f_measure > max_F:
44             max_F = f_measure
45             max_node = node
46             max_rank = restricted_taxonomy.nodes[node].get('rank')
47
48     r = f"The optimal assignment of {read_id
49     } is {max_node} with F_measure = {max_F}, taxonomic rank {max_rank} "
50     results.append(r)
51     optimal_ranks[read_id] = max_rank # Store the rank
52
53 # Print results
54 print("OPTIMAL TAXONOMIC ASSIGNMENTS")
55 for result in results:
56     print(result)
57
58 highest_rank = min(optimal_ranks.values(), key=lambda x: valid_ranks.index(x))
59 print(f"Highest taxonomic rank among optimal assignments: {highest_rank}")
60
61 lowest_rank = max(optimal_ranks.values(), key=lambda x: valid_ranks.index(x))
62 print(f"Lowest taxonomic rank among optimal assignments: {lowest_rank}")

```

Listing 4: Python code for ex 4

Results after running the code are shown in the image below (Figure 3)

```

OPTIMAL TAXONOMIC ASSIGNMENTS
The optimal assignment of R00010 is 7088 with F_measure = 3.289419581914771e-05, taxonomic rank order
The optimal assignment of R00020 is 7100 with F_measure = 0.0005665722379603399, taxonomic rank family
The optimal assignment of R00030 is 7458 with F_measure = 0.0012606366214938543, taxonomic rank family
The optimal assignment of R00040 is 953 with F_measure = 0.0011940298507462687, taxonomic rank genus
The optimal assignment of R00050 is 10193 with F_measure = 0.005641748942172073, taxonomic rank family
The optimal assignment of R00060 is 7100 with F_measure = 0.0005633802816901409, taxonomic rank family
The optimal assignment of R00070 is 384774 with F_measure = 0.01990049751243781, taxonomic rank genus
The optimal assignment of R00080 is 4747 with F_measure = 0.0003652634462606155, taxonomic rank family
The optimal assignment of R00090 is 337673 with F_measure = 0.017699115044247787, taxonomic rank family
The optimal assignment of R00100 is 9397 with F_measure = 0.0018009905447996398, taxonomic rank order
Highest taxonomic rank among optimal assignments: order
Lowest taxonomic rank among optimal assignments: genus

```

Figure 3: Results ex 4

5. **What is the highest taxonomic rank (that is, toward kingdom) for these taxonomic assignments?**  
The highest is order.
6. **What is the lowest taxonomic rank (that is, toward species) for these taxonomic assignments?**  
The lowest is genus.