

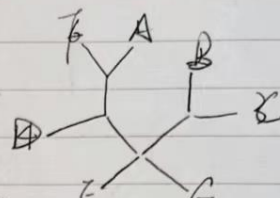
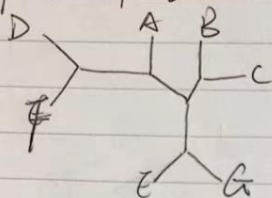
Interpretation of my implementation of Day's algorithm 1985:

Author: Zhang Leyi, Master's in Aarhus University BiRC

2022-04-20

AiB Project 4: Day's Algorithm.

Compute RF Distance for two unrooted trees. In newick format.



$T_1 = (((D, F), A), (B, C), (E, G));$
 $T_2 = (((A, F), D), (B, C), (E, G));$

① Store these two unrooted trees. $O(n)$

Using adjacency list + stack.

def unrootReader(newick)

Read each symbol and letter i from right to left in newick

if $i == ')$ # should push in.

using na to store the strings I read.

if $na == ""$, only push $)$ into S_1 .

else push $)$ into S_1 and push na into S_2 .

elif $i == ','$ # A, B

same as $)$

elif $i == '('$

push in

pop until $)$ # pop both S_1 and S_2 . and all strings popped are in the same layer (A, B) \Rightarrow

else

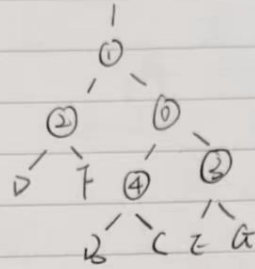
na = i + na.

In doing elif $i == '('$, we can put (A, B)C into adjacency
 $A \rightarrow C$, $C \rightarrow A$, $B \rightarrow C$, $C \rightarrow B$

In order to DFS later.

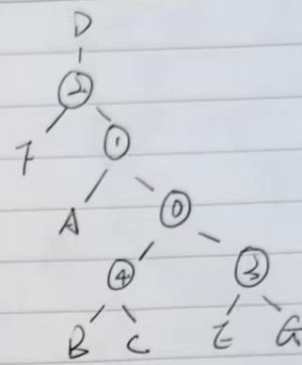
② Since two trees should have same names in leaves, select one leaf as root and DFS to transfer them into rooted tree.
 $O(n)$.

After DFS $T_1: A$



Just DFS.

$T_2:$

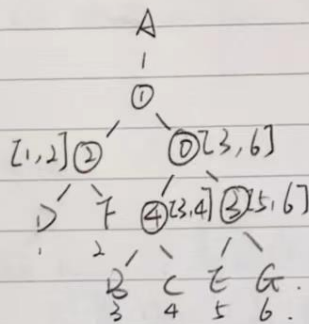


③ Relabel new leaves by using numbers. $O(n)$

D	F	B	C	E	G
1	2	3	4	5	6

still DFS.

④ Calculate the scales $O(n)$



Using dies to make map on ~~leaves~~ ^{nodes and} number.
Here it used minscale, maxscale, size.

DFS to search for every node u and find father.
 $\text{minscale}[father] = \min(\text{minscale}[u], \text{minscale}[f])$
 maxscale —

See line 40

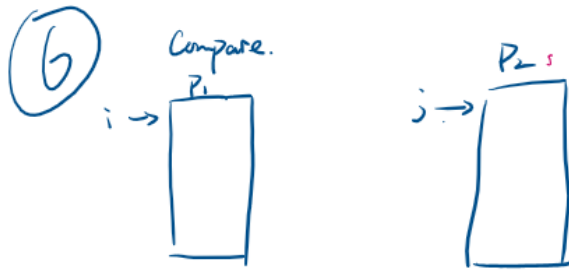
⑤ Radix sort these scales. $O(n)$

Right now we have scales of T_1 and T_2

$\Rightarrow [a, b]$: a is the first key, ~~to~~ b is the second key.

[1,2] ↓
[3,4] ↓
[3,6] ↓
[5,6] ↓

⑥ Find shares $RF = |T_1| + |T_2| - 2 \cdot \text{shares}$ $O(n)$.



if $P_1(i) > P_2(j)$

$j++$ # $P_1(i)$ could possibly exist in P_2 later.

elif $P_1(i) == P_2(j)$

share += 1.

$i++$; $j++$.

else # $P_1(i) < P_2(j)$ # could not be existed

$i++$.

Until $i == \text{len}(P_1)$ or $j == \text{len}(P_2)$

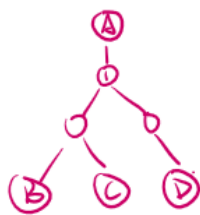
return share.

⑦ $|T_1| + |T_2| - 2 \cdot \text{share}$

$|T_1|$ and $|T_2|$ are internal nodes.

$T_1 = \text{len}(\text{nodes}) - \text{len}(\text{leaves}) - 2$

Why minus 2:



A and C are not seen as internal nodes here.

RT distance = $T_1 + T_2 - 2 \times \text{shares}$