

Thoughts on constructing the Full Supertree

October 23, 2015

In the *otcetera* pipeline, the synthesis tree is first constructed using a pruned taxonomy. This taxonomy has been pruned to remove any leaves that do not occur in one of the ranked input phylogenies. Therefore, a later step in the pipeline involves adding the pruned taxa back in to the synthesis tree. The resulting tree is the *full supertree*. This document is an attempt to collect, and perhaps organize, thoughts on how the full supertree can or should be constructed, as well as related questions and concepts.

This document takes the approach that we should be able to use the sub-problem solver to construct the full supertree by feeding it a sequence of 2 trees:

```
otc-solve-subproblem grafted-solution.tre cleaned_ott.tre
```

Various other quick-and-dirty methods may be sufficient to achieve the result. However, the attempt to make the subproblem solver fast enough to solve this particular problem raises various issues with the subproblem solver. It also suggests various optimizations that may be useful more generally.

1 Potential speed increases

In the current implementation of the BUILD algorithm, we have a number of places we take excessive computation time:

- We attempt to construct rooted splits (a.k.a. `desIds`) for each node. For a bifurcating tree, this operation should take time and memory quadratic in the number of leaves.
- We attempt to construct connected component by considering each split in a tree separately. However, considering splits for nodes that are not direct children of the root is redundant.
- Much of the time is spent in determining which splits are imposed at a given level in the tree, and therefore need not be passed to subproblems.
 - When different splits have the same leaf set, we should be able to get a speedup.
 - When some splits have the full leaf set, we should be able to get a speedup.
- We recompute connected components from scratch each time BUILD is recursively called on a subproblem. This could be avoided by incrementally removing edges from the graph and discovering new connected components that appear, as in Henzinger et al. *However, it is unclear if an algorithm similar to Henzinger et al could be used to find the edges to remove at each step.*
- When we have two trees T_1 and T_2 , and either $\mathcal{L}(T_1) \subseteq \mathcal{L}(T_2)$ or $\mathcal{L}(T_2) \subseteq \mathcal{L}(T_1)$, then it should be possible to determine all conflicting splits in a single pass over the trees, similar to `otc-detectcontested`.

2 The problem

When the subproblem to be solved consists of two ranked trees, T_1 and T_2 , and the second tree is the taxonomy, then we have $\mathcal{L}(T_1) \subseteq \mathcal{L}(T_2)$. For each (rooted) split in T_2 , we can determine whether that split is consistent with T_1 . We claim that each split in T_2 is either consistent with T_1 , or incompatible with at least one split of T_1 . We can therefore remove each split of T_2 that is inconsistent with T_1 to form a new tree T'_2 . The splits of T_1 and T'_2 are then jointly consistent. Furthermore, by combining the trees T_1 and T'_2 , we obtain a new tree in which the splits of T_1 that are not implied by T'_2 may not fully specify where certain taxa of T'_2 are placed. We resolve this ambiguity by placing such taxa rootward, but their range of attachment extends over specific branches in the tree obtained by combining T_1 and T'_2 . All of these branches derive from T_1 but not from T'_2 .

3 Questions

- Is there a single, unique solution to this problem?
- When running BUILD, is it possible to construct a *reason* for the lack of inclusion of each split? Specifically, can we say which split (or set of splits) conflicts with that split?