

2.2 Outliers for multiple related taxa

First, generate a list of square distance matrices.

```
> dmats <- lapply(taxa, function(taxon) {
+   ape::dist.dna(seqs[taxon$i,], pairwise.deletion=TRUE, as.matrix=TRUE, model='raw')
+ })
```

Calculate outliers for each matrix. Here (as above) we are using a distance threshold of 1.5% from the “center-most” sequence (i.e., the one with the least sum of pairwise distances to every other sequence).

```
> outliers <- sapply(dmats, findOutliers, cutoff=0.015)
```

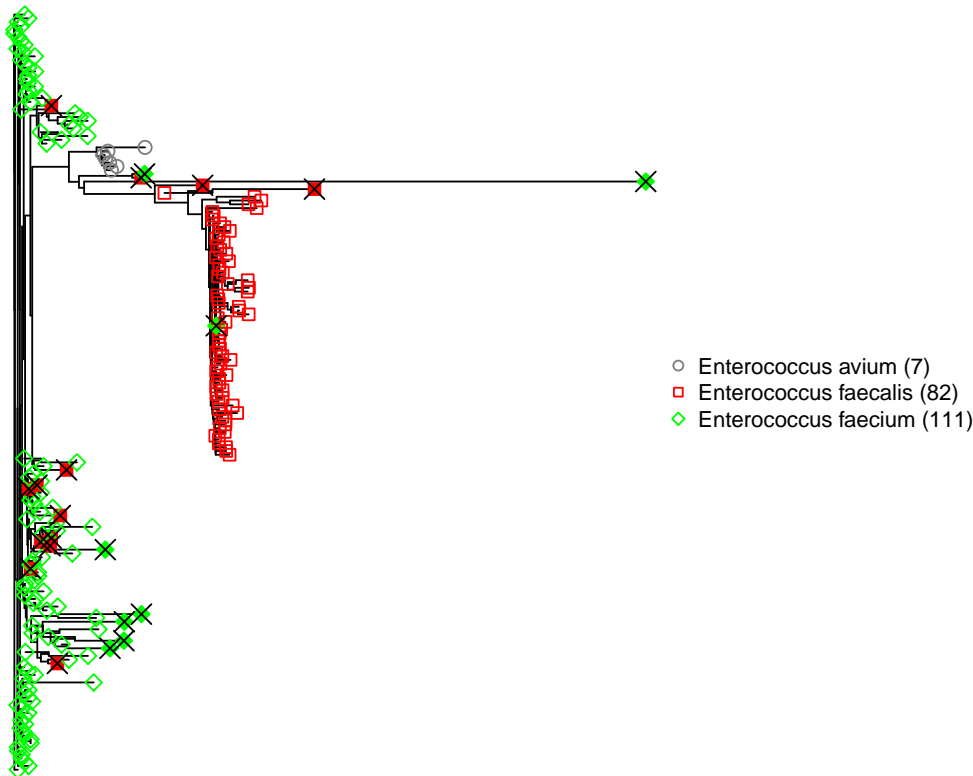
Add results of outlier status to seqdat.

```
> seqdat$outlier <- FALSE
> for(x in outliers){
+   seqdat[match(names(x), seqdat$seqname), 'outlier'] <- x
+ }
> with(seqdat, table(tax_name, outlier))
```

	outlier	
tax_name	FALSE	TRUE
Enterococcus avium	7	0

Finally, we can visualize the fact that many of the outliers are actually the result of labels being switched between taxa (that is, *E. faecium* sequences are labeled as *E. faecalis*) and vice versa. In the tree below, terminal nodes are identified according to the original species labels.

```
> with(seqdat, {
+   dmat <- ape::dist.dna(seqs, pairwise.deletion=TRUE, as.matrix=TRUE, model='raw')
+   clstutils::prettyTree(nj(dmat), groups=tax_name,
+                         ## type='unrooted',
+                         X=outlier, fill=outlier)
+ })
>
```



3 Selecting a diverse subset

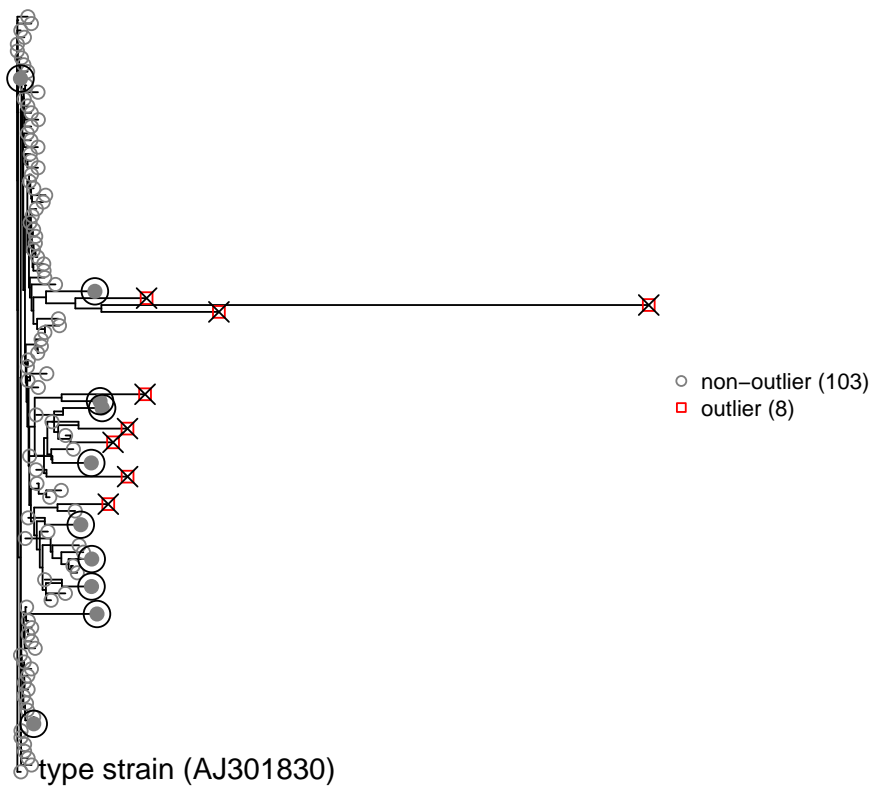
Because we cannot use every available sequence in our reference tree, a sampling strategy is required. One strategy is to select a maximally diverse subset of sequences. The function `clstutils::maxDists` performs this operation. In addition, we can exclude sequences identified as outliers in the previous step (outlier identification is critical here, lest we select primarily outliers!). We can also optionally include the “centermost” sequence in the set, plus any type strains.

```
> with(seqdat[Efaecium,], {
+   selected <- clstutils::maxDists(dmat, idx=which(isType),
```

```

+                                     N=10, exclude=outlier, include.center=TRUE)
+ prettyTree(nj(dmat), groups=ifelse(outlier,'outlier','non-outlier'),
+         X=outlier,
+         O=selected, fill=selected,
+         labels=ifelse(isType,gettextf('type strain (%s)', accession),NA))
+ })

```



Here the selected sequences are identified with circled, filled glyphs.