

Impact of tree choice in metagenomics differential abundance studies

Antoine Bichat^{1, 2, 3}, Christophe Ambroise², Jonathan Plassais¹,
and Mahendra Mariadassou³

¹Entérome, Paris, France

²LaMME, Université Evry Val d'Essonne, Evry, France

³INRA - MaIAGE, Université Paris Saclay, Jouy-en-Josas, France

We consider the problem of incorporating evolutionary information (e.g. taxonomic or phylogenetic trees) in the context of metagenomics differential analysis. Recent results published in the literature propose different ways to leverage the tree structure to increase the detection rate of differentially abundant taxa. Here, we propose instead to use a different hierarchical structure, in the form of a correlation-based tree, as it may capture the structure of the data better than the phylogeny. We first show that the correlation tree and the phylogeny are significantly different before turning to the impact of tree choice on detection rates. Using synthetic data, we show that the tree does have an impact: smoothing p-values according to the phylogeny leads to equal or inferior rates as smoothing according to the correlation tree. However, both trees are outperformed by the classical, non hierarchical, Benjamini-Hochberg (BH) procedure in terms of detection rates. Other procedures may use the hierarchical structure with profit but do not control the False Discovery Rate (FDR) a priori and remain inferior to a classical Benjamini-Hochberg procedure with the same nominal FDR. On real datasets, no hierarchical procedure had significantly higher detection rate than BH. Although intuition advocates the use of a hierarchical structure, be it the phylogeny or the correlation tree, to increase the detection rate in microbiome studies, current hierarchical procedures are still inferior to non hierarchical ones and effective procedures remain to be invented.