

Using the concept of informative genomic segment to investigate microbial diversity of metagenomics sample

December 18, 2014

Abstract

In almost all the metagenomics projects, diversity analysis plays an important role to supply information about the richness of species, the species abundance distribution in a sample or the similarity and difference between different samples, all of which are crucial to draw insightful and reliable conclusion. Traditionally especially for amplicon metagenomics data set, OTUs(Operational Taxonomic Units) based on 16S rRNA genes are used as the cornerstone for diversity analysis. Here we propose a novel concept - IGS (informative genomic segment) and use IGS as a replacement of OTUs to be the cornerstone for diversity analysis of whole shotgun metagenomics data sets. IGSs represent the unique information in a metagenomics data set and the abundance of IGSs in different samples can be retrieved by the reads coverage through an efficient k-mer counting method. This samples-by-IGS abundance data matrix is a promising replacement of samples-by-OTU data matrix used in 16S rRNA based analysis and all existing statistical methods can be borrowed to work on the samples-by-IGS data matrix to investigate the diversity. We applied the IGS-based method to several simulated data sets and a real data set - Global Ocean Sampling Expedition (GOS) to do beta-diversity analysis and the samples were clustered more accurately than existing alignment-based method. We also tried this novel method to Great Prairie Soil Metagenome Grand Challenge data sets. Furthermore we will show some preliminary results using the IGS-based method for alpha-diversity analysis. Since this method is totally binning-free, assembly-free, annotation-free, reference-free, it is specifically promising to deal with the highly diverse samples, while we are facing large amount of dark matters in it, like soil.

1 Introduction

2 Results

2.1 IGS(informative genomic segment) can represent the novel information of a genome

2.2 IGS can be used to do alpha diversity analysis

2.3 IGS can be used to do beta diversity analysis

2.4 GOS data sets: Sorcerer II Global Ocean Sampling Expedition

2.5 MetHit Human Gut metagenomics data set

2.6 HMP metagenomics data set

2.7 GPGC and ARMO soil sample

3 Discussion

3.1 IGS can be the foundation of a new framework to do microbial diversity analysis

3.2 IGS is promising to more problems

3.3 Concluding thoughts

4 Conclusion

5 Methods

6 Acknowledgments