

Mini-metagenomes assembled by SPAdes

Margarita Akseshina

scientific advisor: Anton Korobeynikov

Goals

- find features for binning
- check them for reliability and robustness
- try to improve metagenomic assembly

Features

- tetranucleotide frequencies
 - 136 dimensions
- average coverage
- contig length
- codon usage
- ...

Datasets

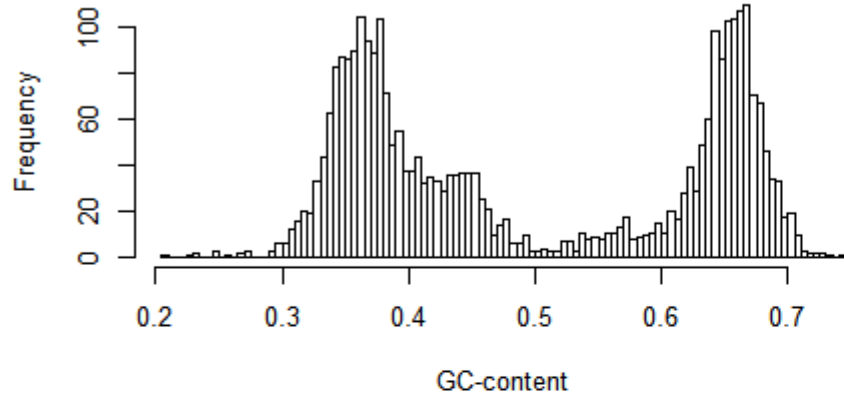
Mini-metagenomes from Lena Gerwick:

- Bastimolide producer
- Moorea Producens JHB
- Scytonema Hofmanii (UTEX B2349)
- Moorea Bouillonii PNG

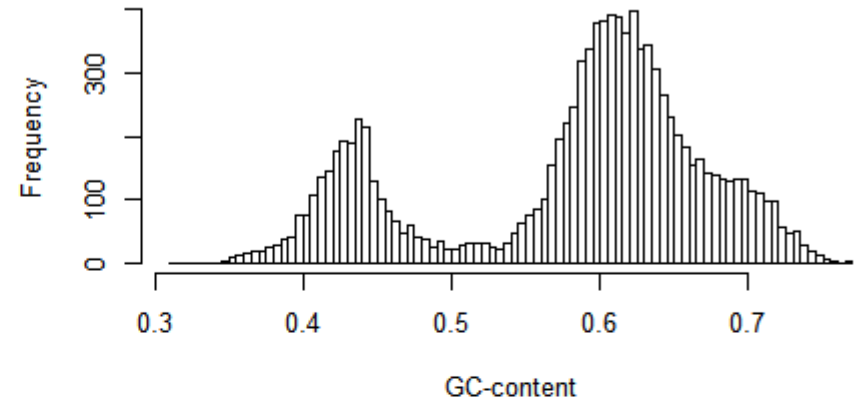
Usually contains a cyanobacteria of interest and several symbionts.

How many taxons do you see?

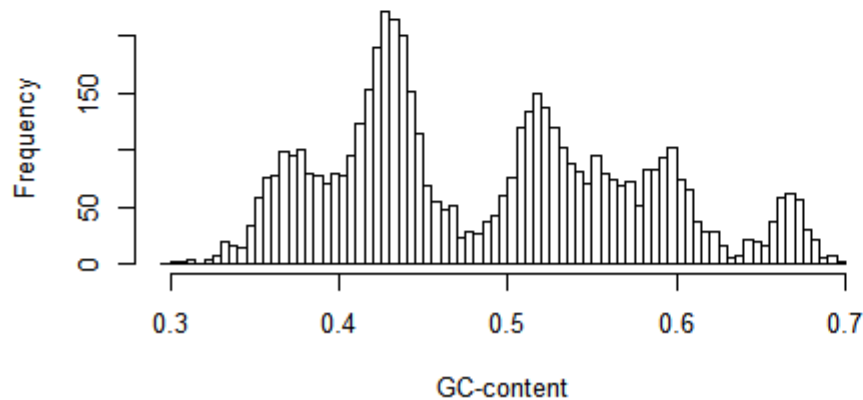
bastimolide



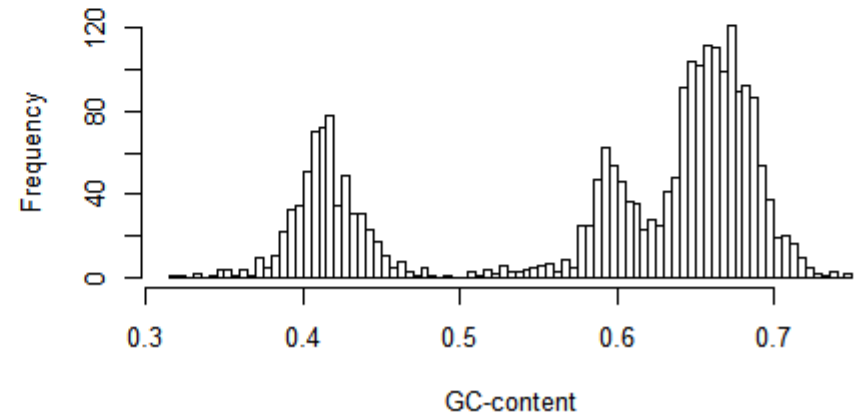
jhb



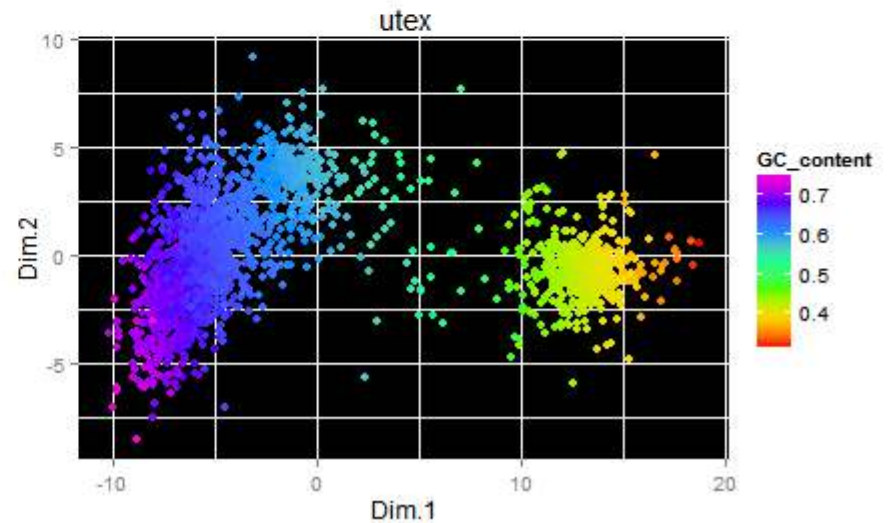
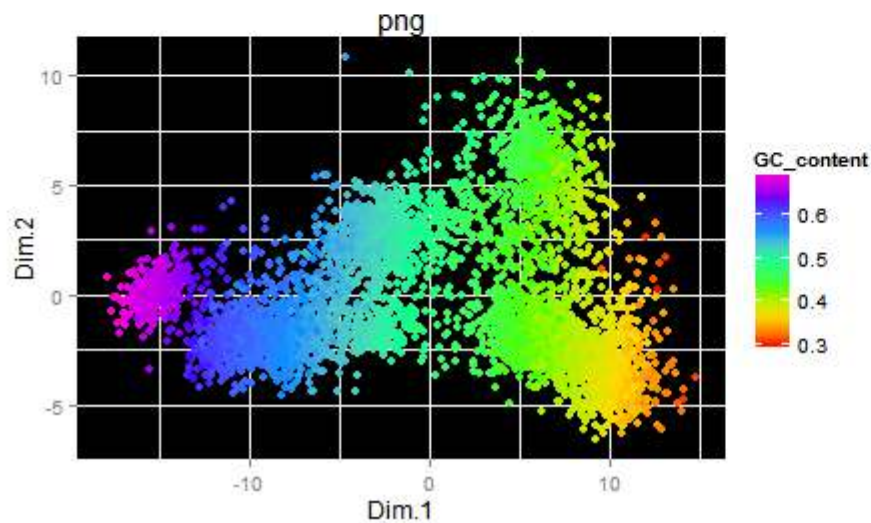
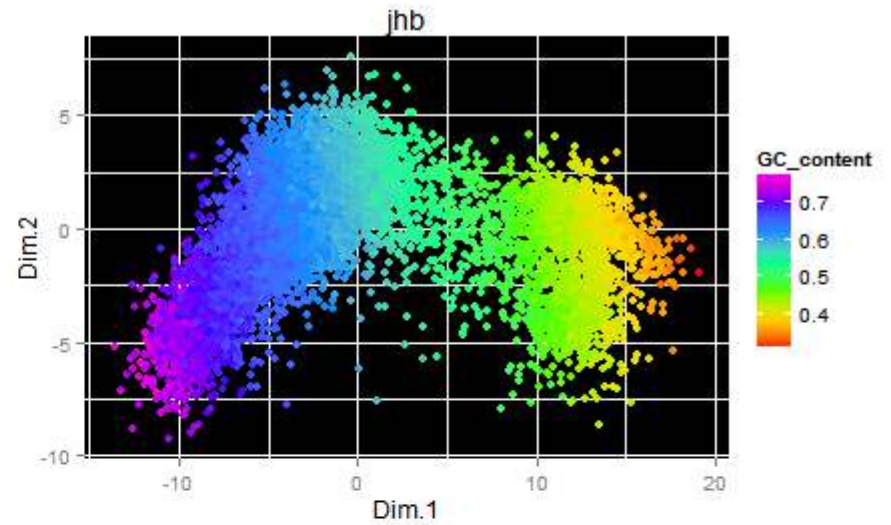
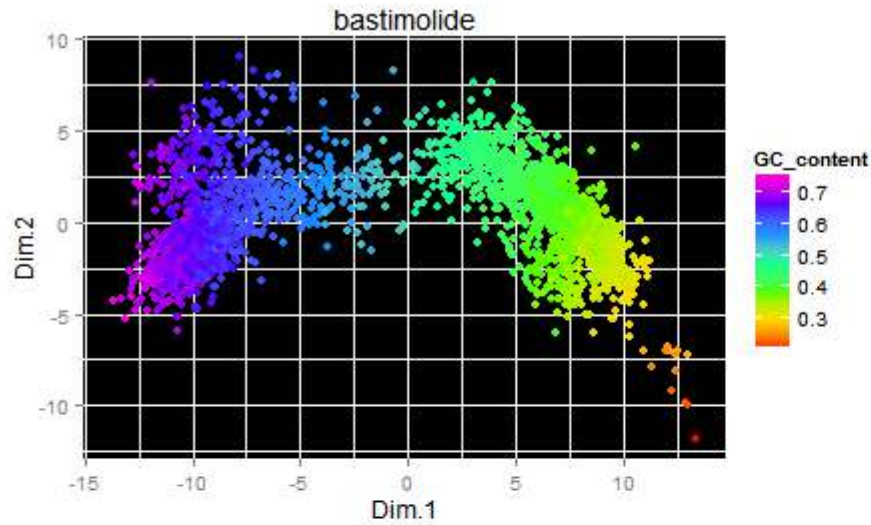
png



utex

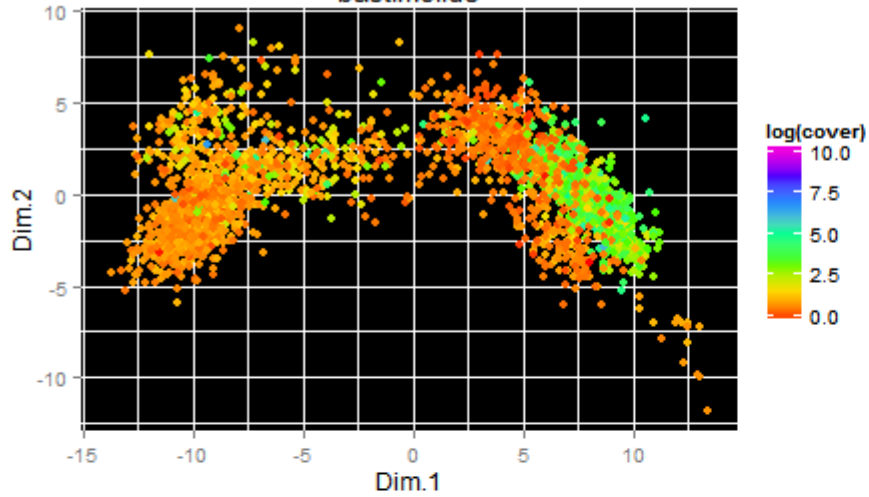


PCA with GC-content

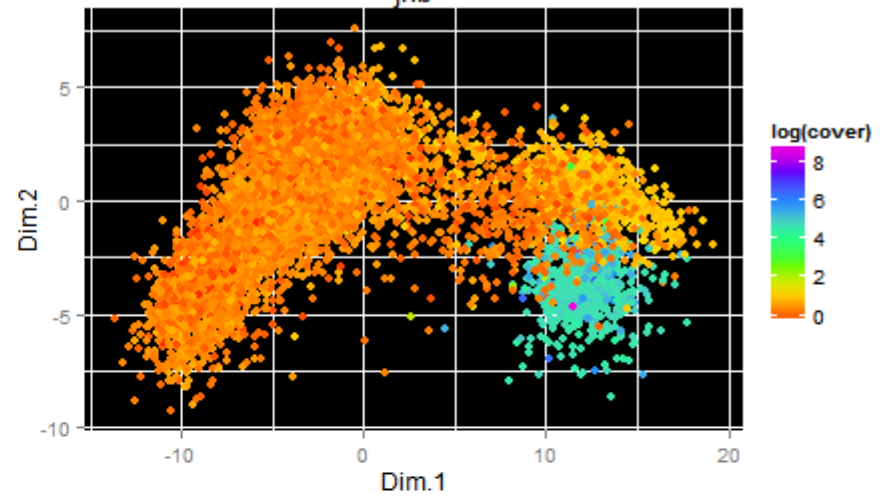


PCA with coverage

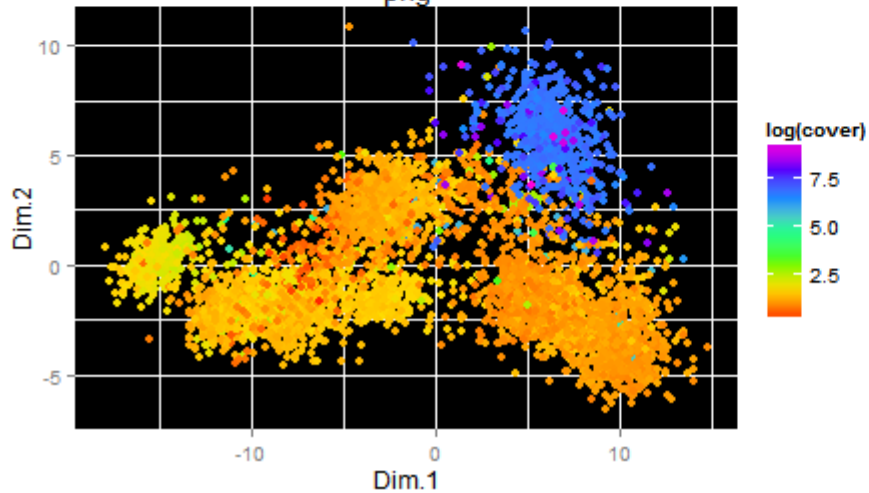
bastimolide



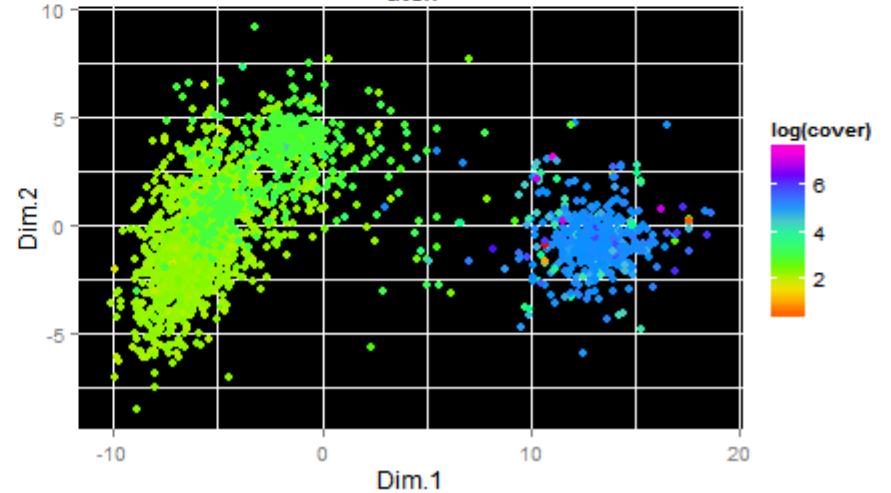
jhb



png



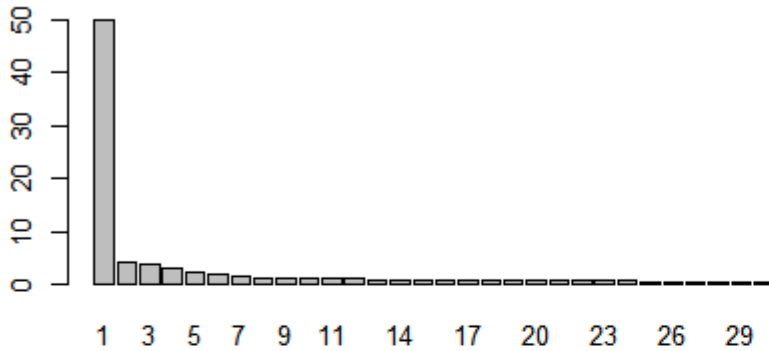
utex



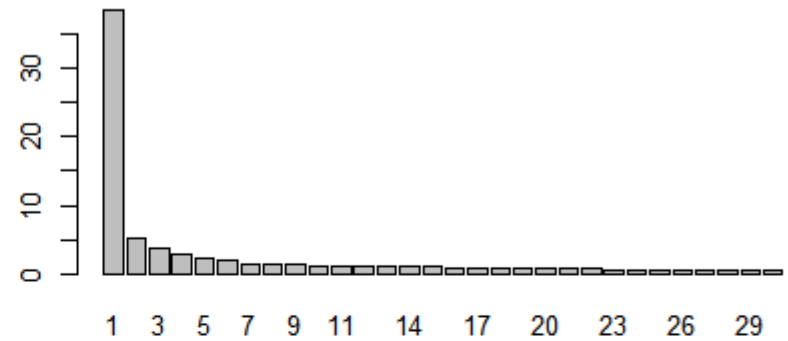
How many PC to keep?

Scree plots

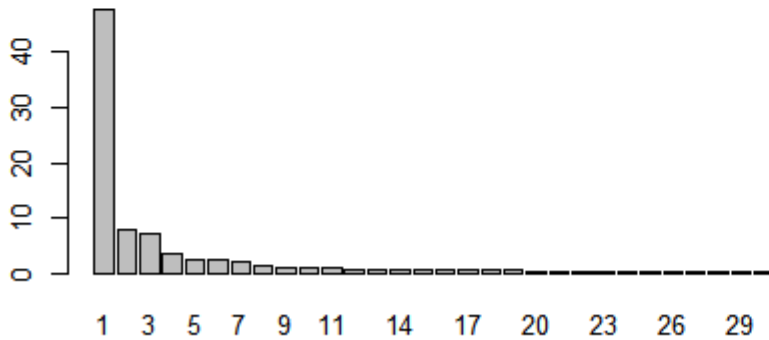
bastimolide



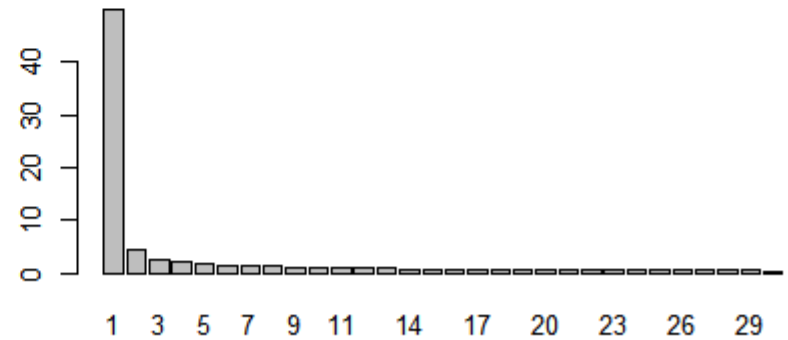
jhb



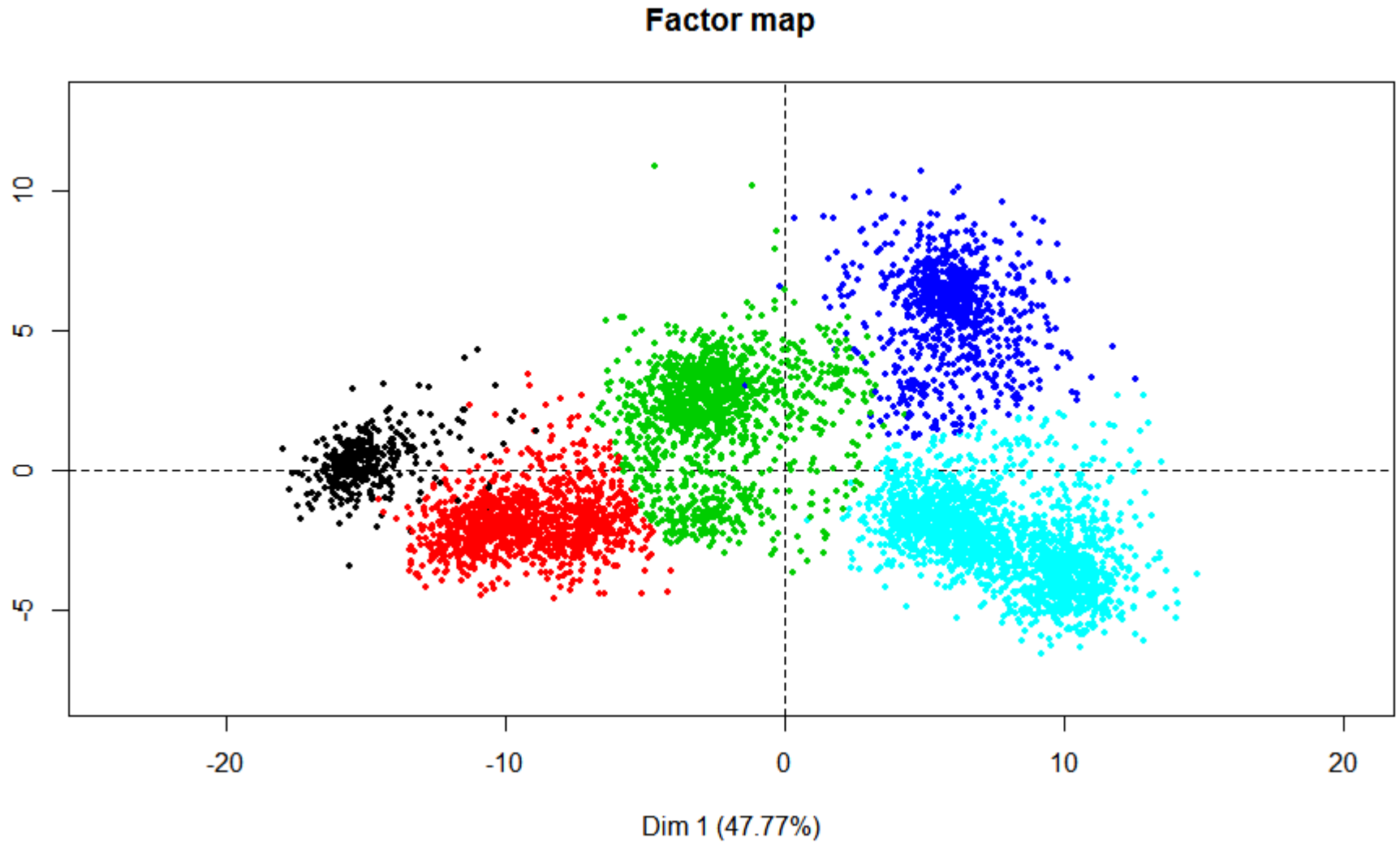
png



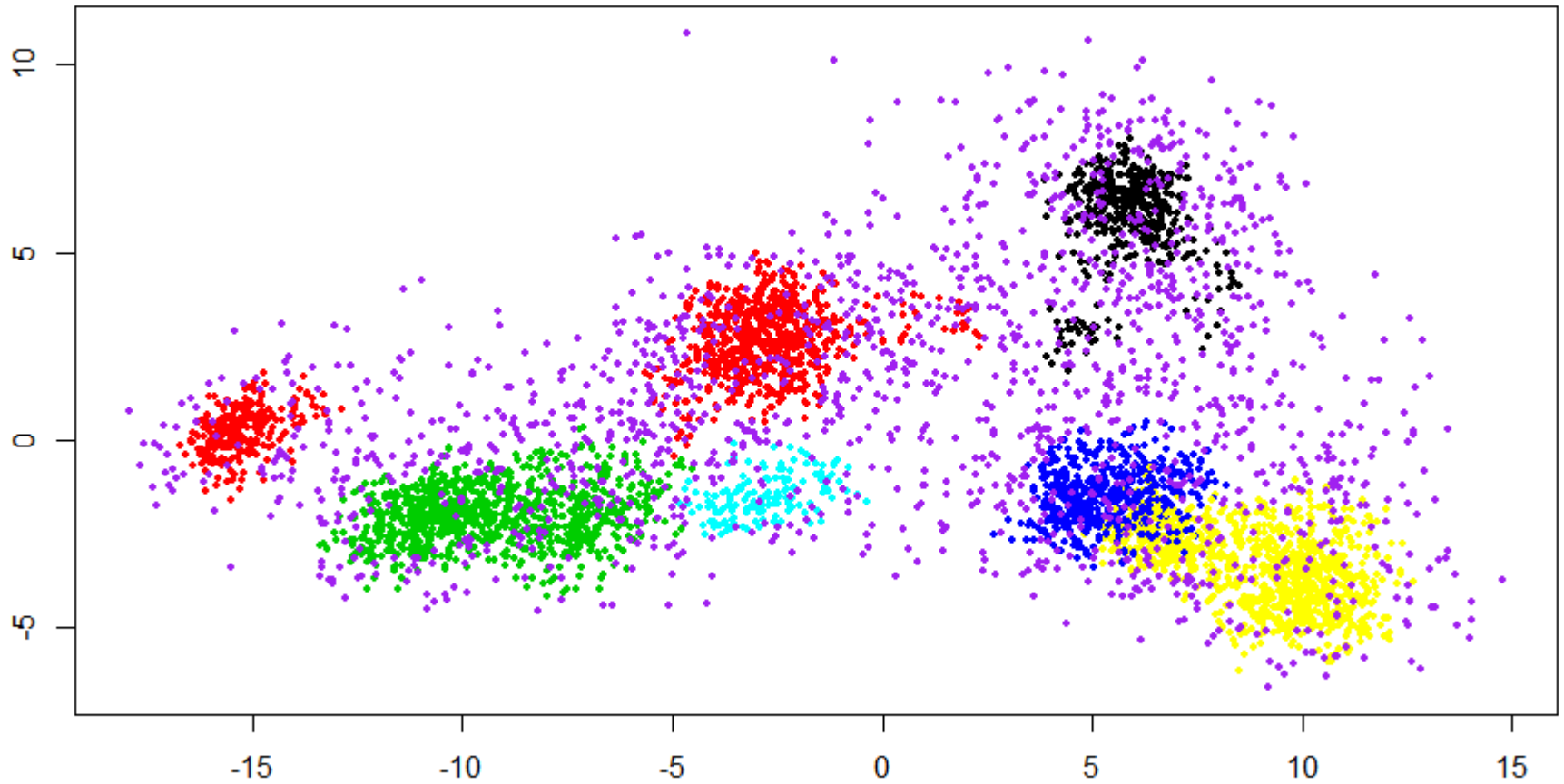
utex



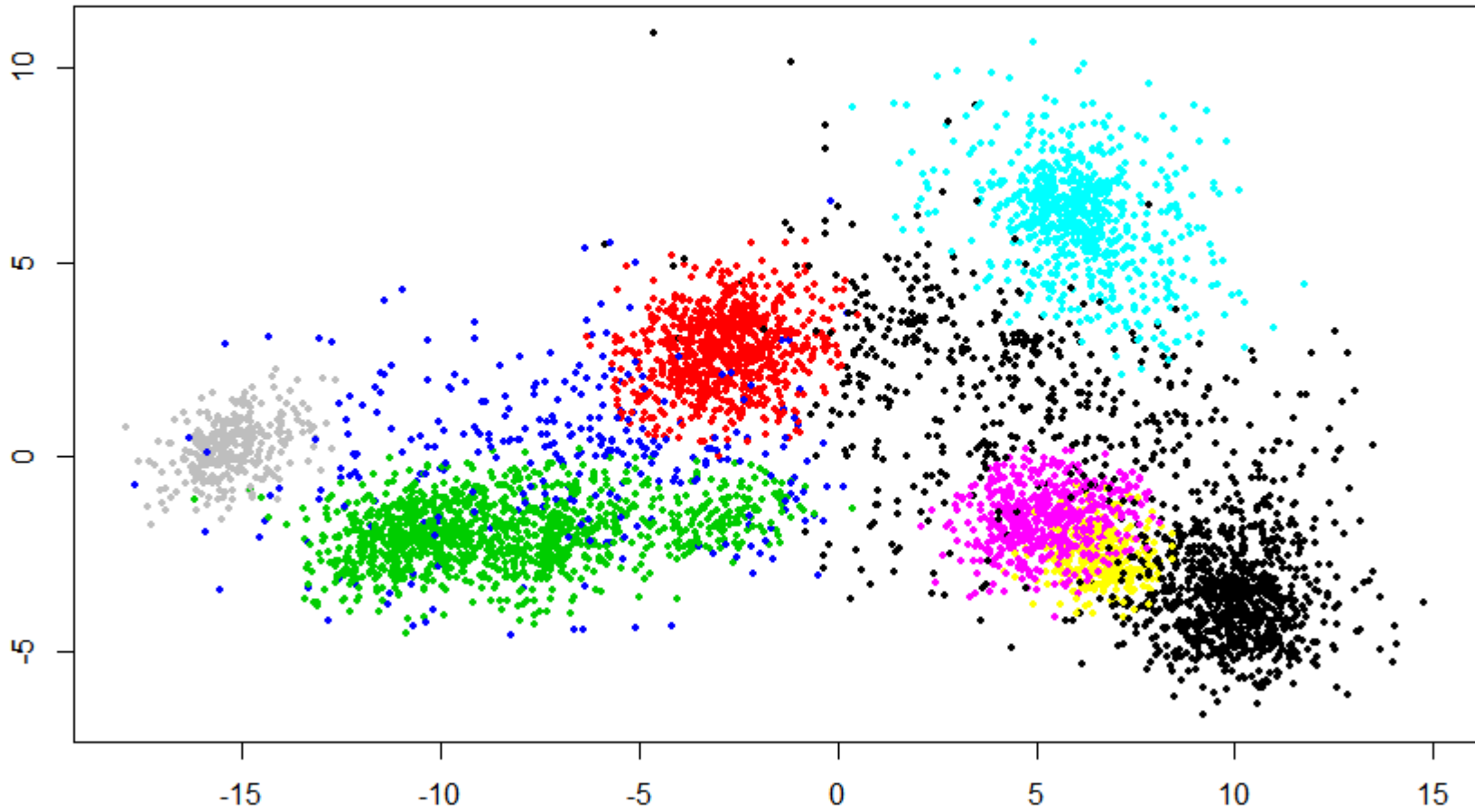
Hierarchical clustering



DBSCAN



Mclust



Further questions

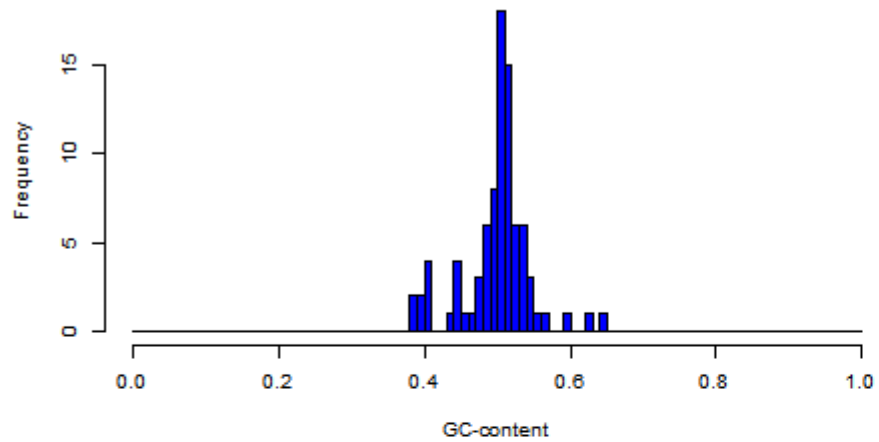
- a suitable length filter
- a method for clustering
- number of principal components
- validation of the results
- ...

How to model data

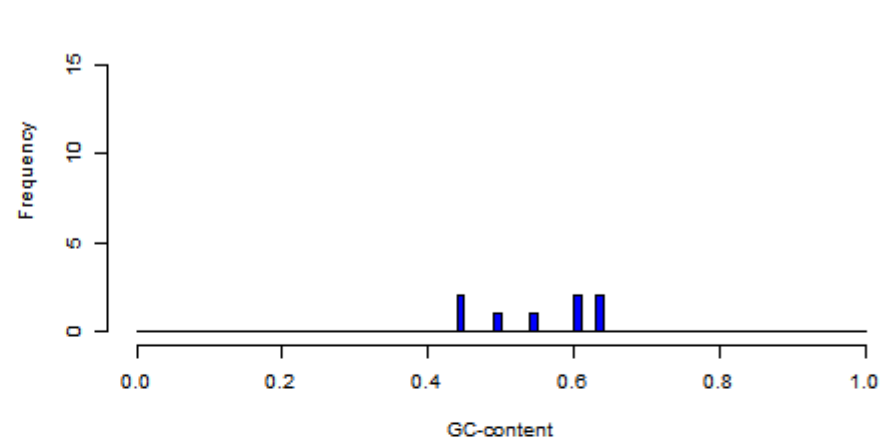
- merge several assemblies
- merge different reads and then assemble

Assemblies

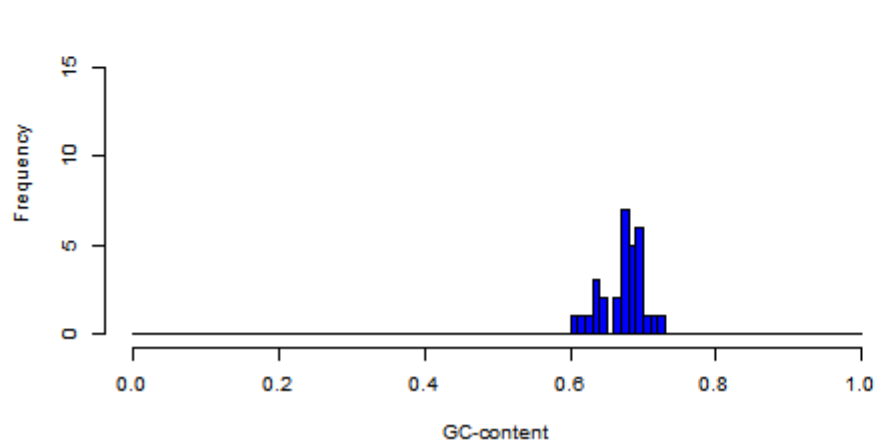
E.coli



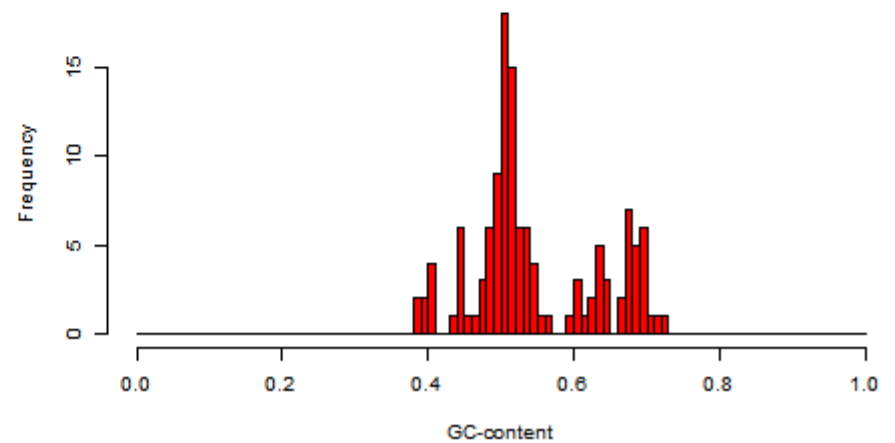
M.ruber



R.sphaeroides

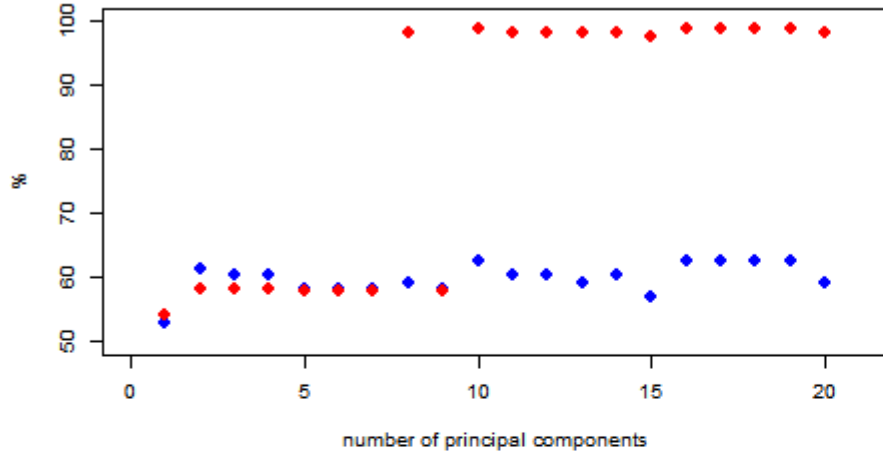


all merged

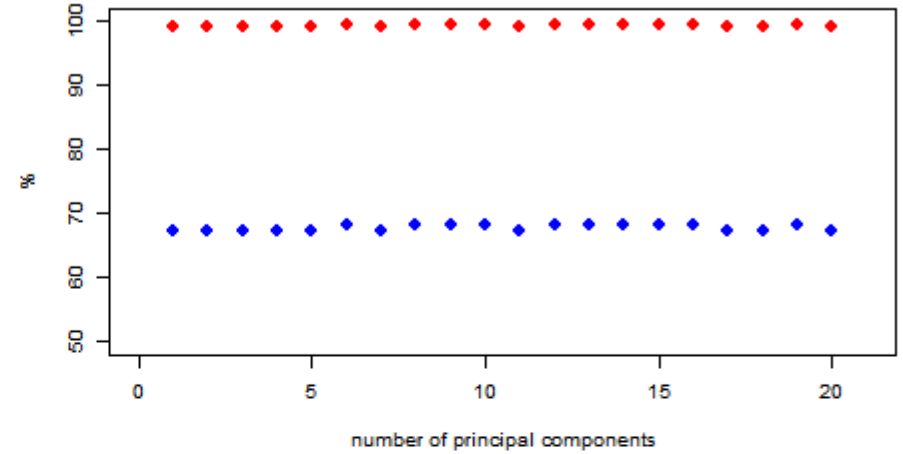


Number of PC (hclust)

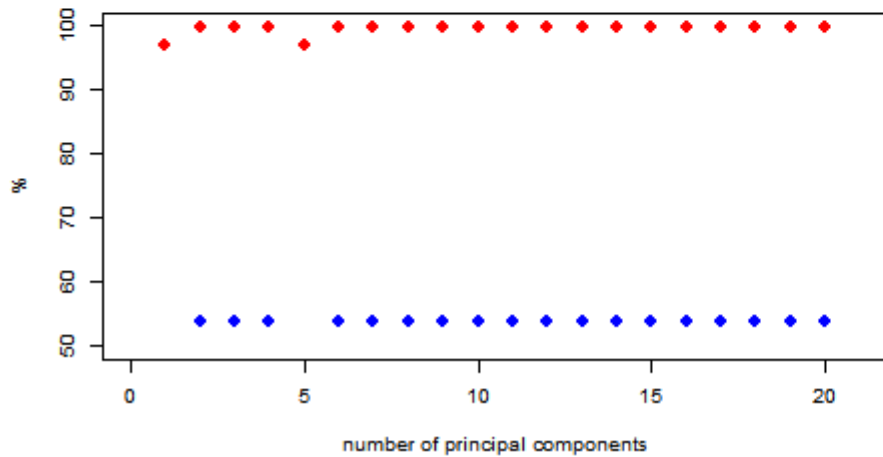
E.coli + M.ruber



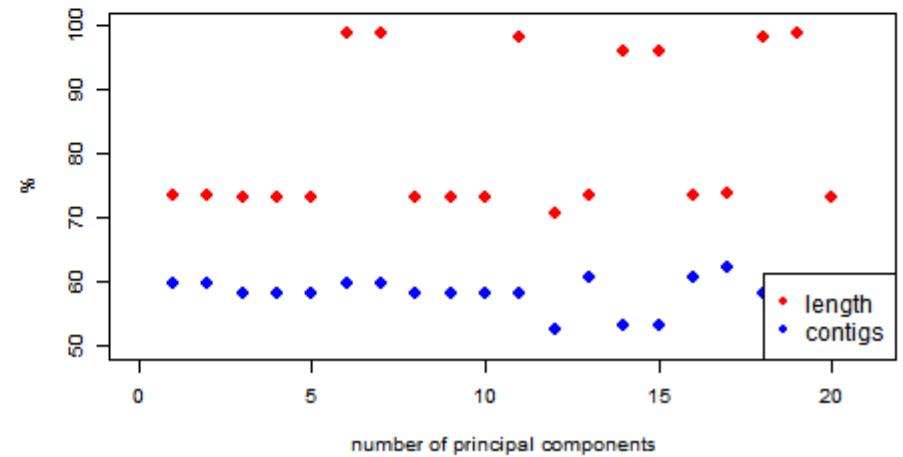
E.coli + R.sphaeroides



M.ruber + R.sphaeroides

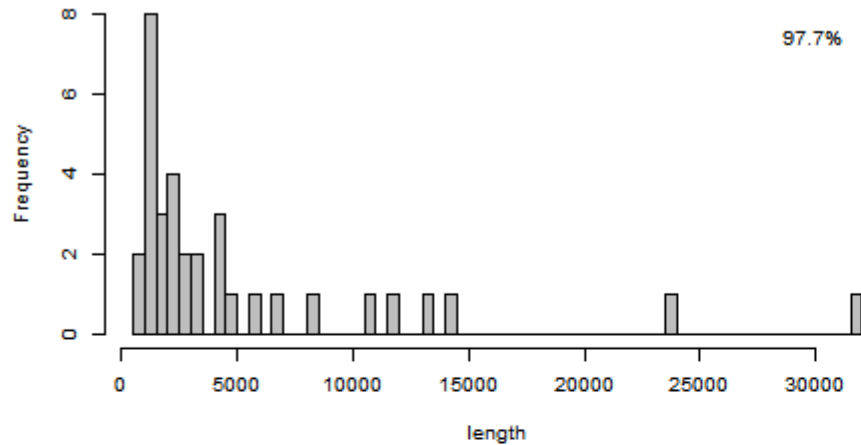


All

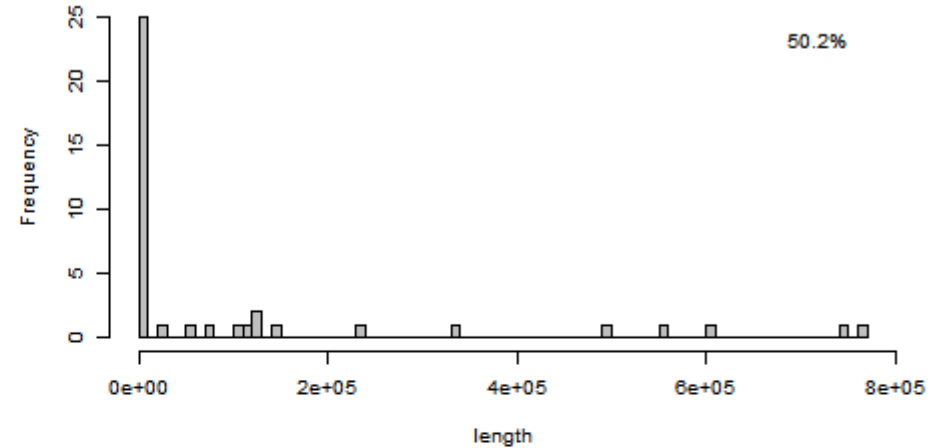


False contigs length (mclust, 10 PC)

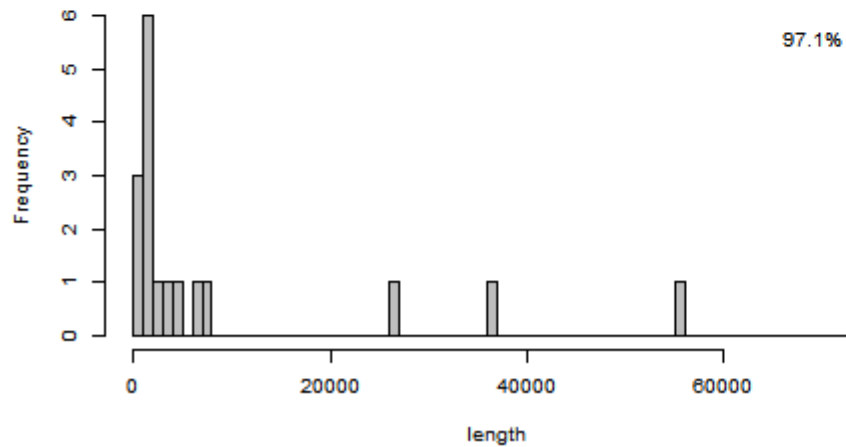
E.coli + M.ruber



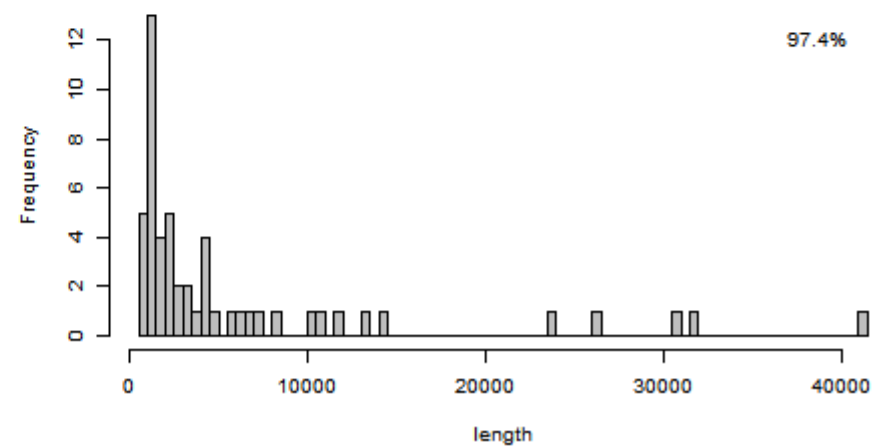
E.coli + R.sphaeroides



M.ruber + R.sphaeroides



All



Thank you!