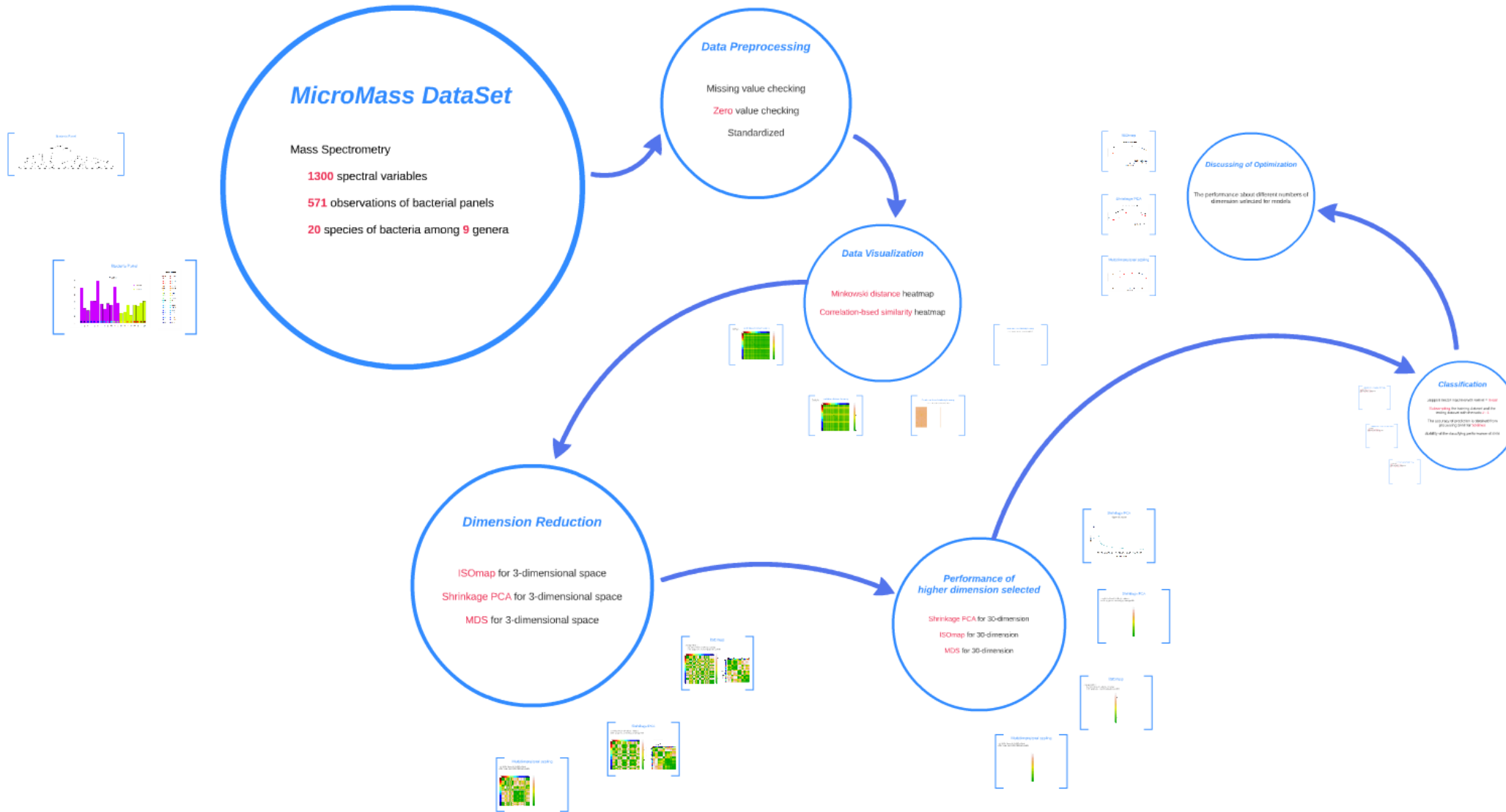


High Dimensional Data Analysis Final Project

THE END
THANK YOU FOR LISTENING




MicroMass DataSet

Mass Spectrometry

1300 spectral variables

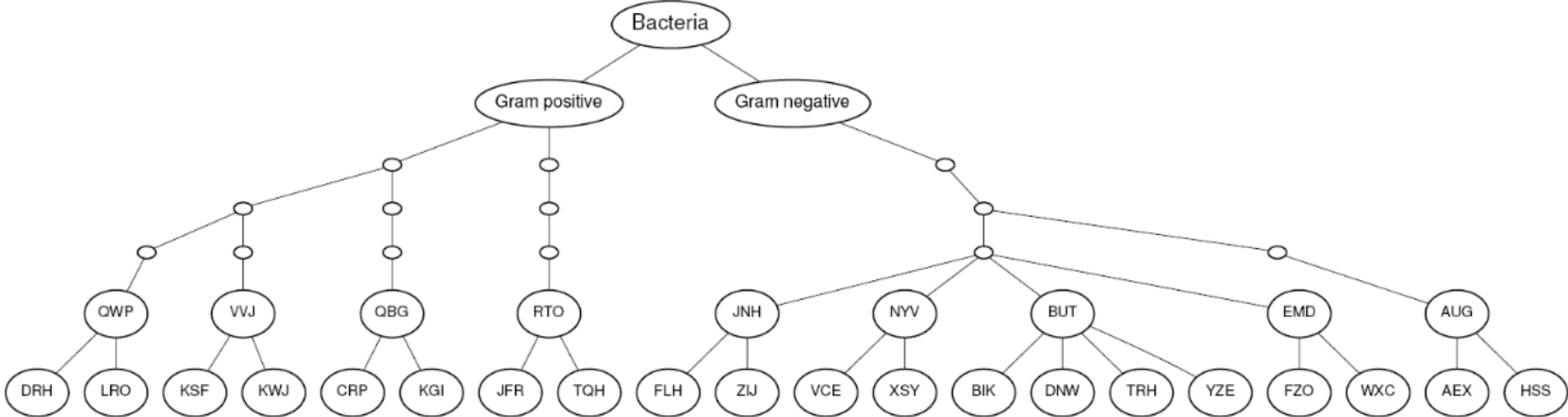
571 observations of bacterial panels

20 species of bacteria among **9** genera

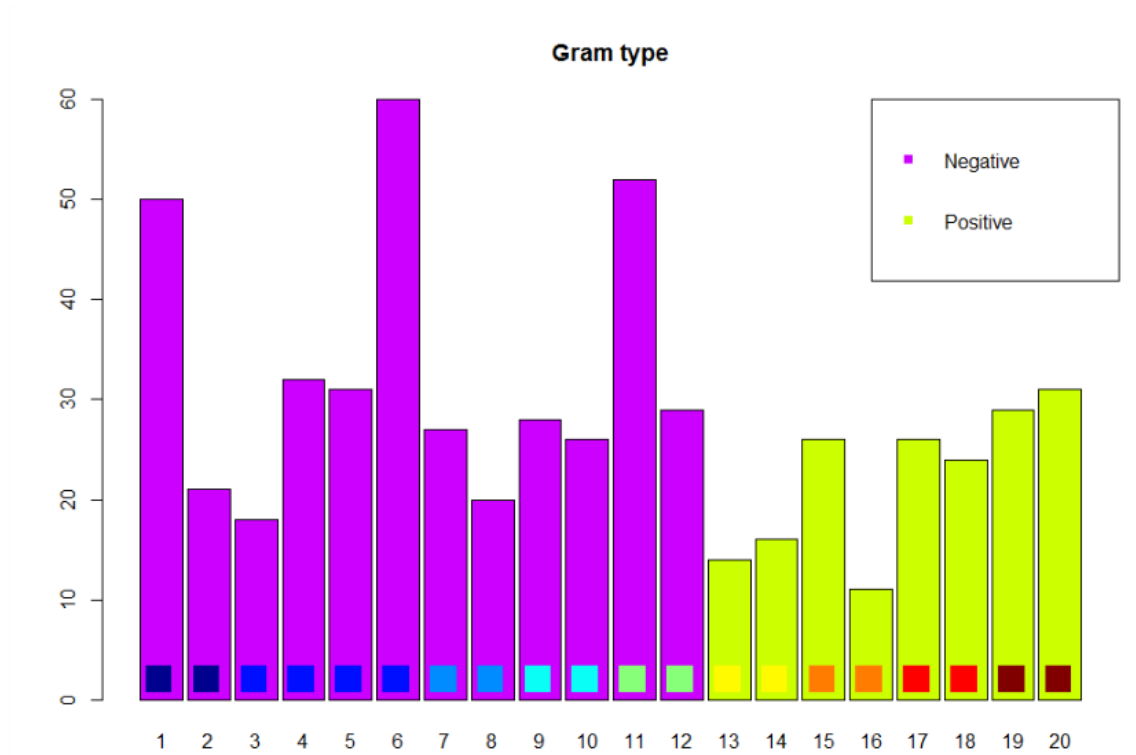


Genus	Species
11-001	11-001
11-002	11-002
11-003	11-003
11-004	11-004
11-005	11-005
11-006	11-006
11-007	11-007
11-008	11-008
11-009	11-009
11-010	11-010
11-011	11-011
11-012	11-012
11-013	11-013
11-014	11-014
11-015	11-015
11-016	11-016
11-017	11-017
11-018	11-018
11-019	11-019
11-020	11-020
11-021	11-021

Bacteria Panel



Bacteria Panel



Genera color table

20	VVJ	VVJ.KWJ
19	VVJ	VVJ.KSF
18	RTO	RTO.TQH
17	RTO	RTO.JFR
16	QWP	QWP.LRO
15	QWP	QWP.DRH
14	QBG	QBG.KGI
13	QBG	QBG.CRP
12	NYV	NYV.XSY
11	NYV	NYV.VCE
10	JNH	JNH.ZJ
9	JNH	JNH.FLH
8	EMD	EMD.WXC
7	EMD	EMD.FZO
6	BUT	BUT.YZE
5	BUT	BUT.TRH
4	BUT	BUT.DNW
3	BUT	BUT.BIK
2	AUG	AUG.HSS
1	AUG	AUG.AEX

Genera Species

Data Preprocessing

Missing value checking

Zero value checking

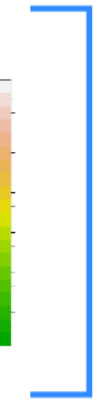
Standardized



Data Visualization

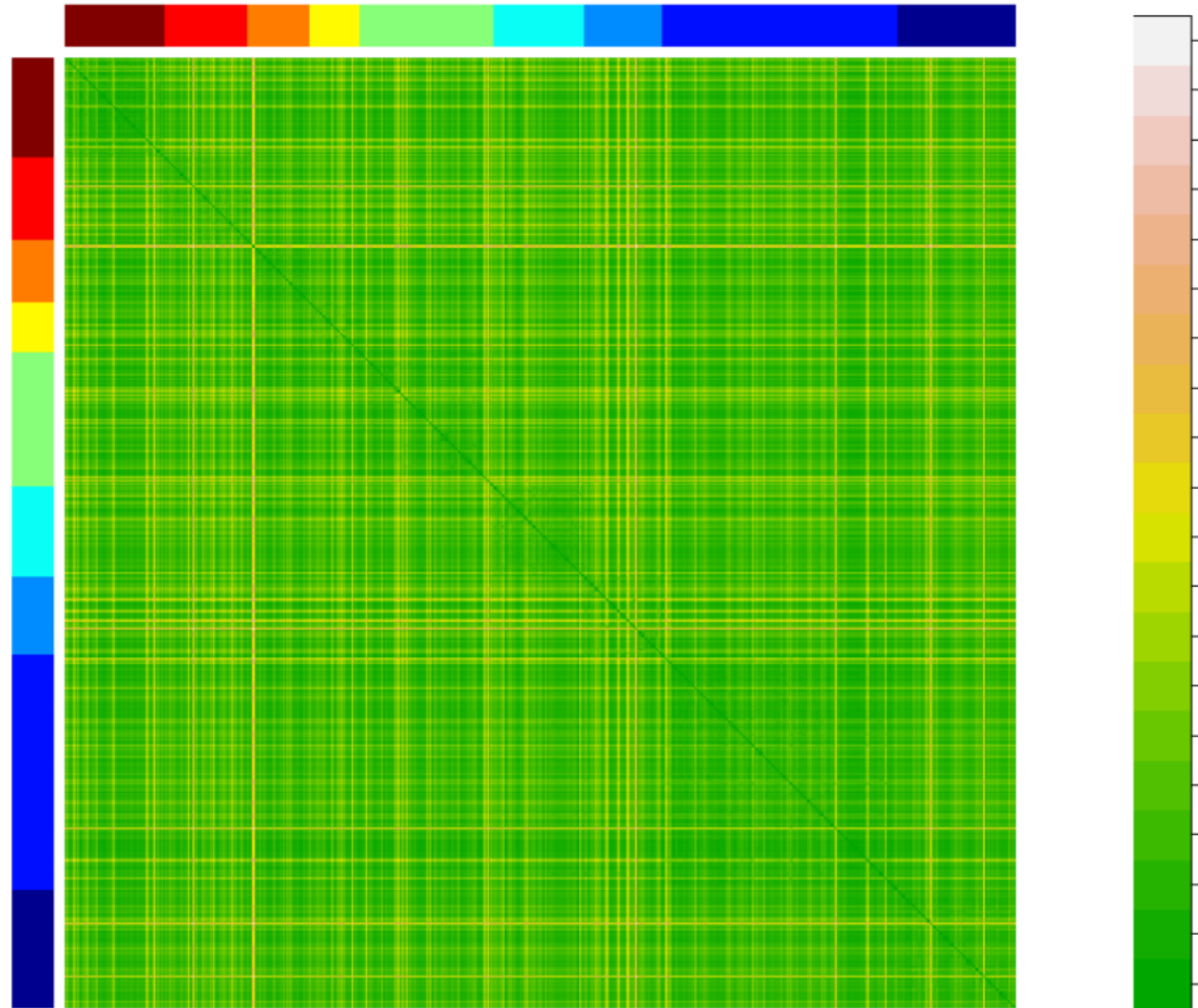
Minkowski distance heatmap

Correlation-based similarity heatmap



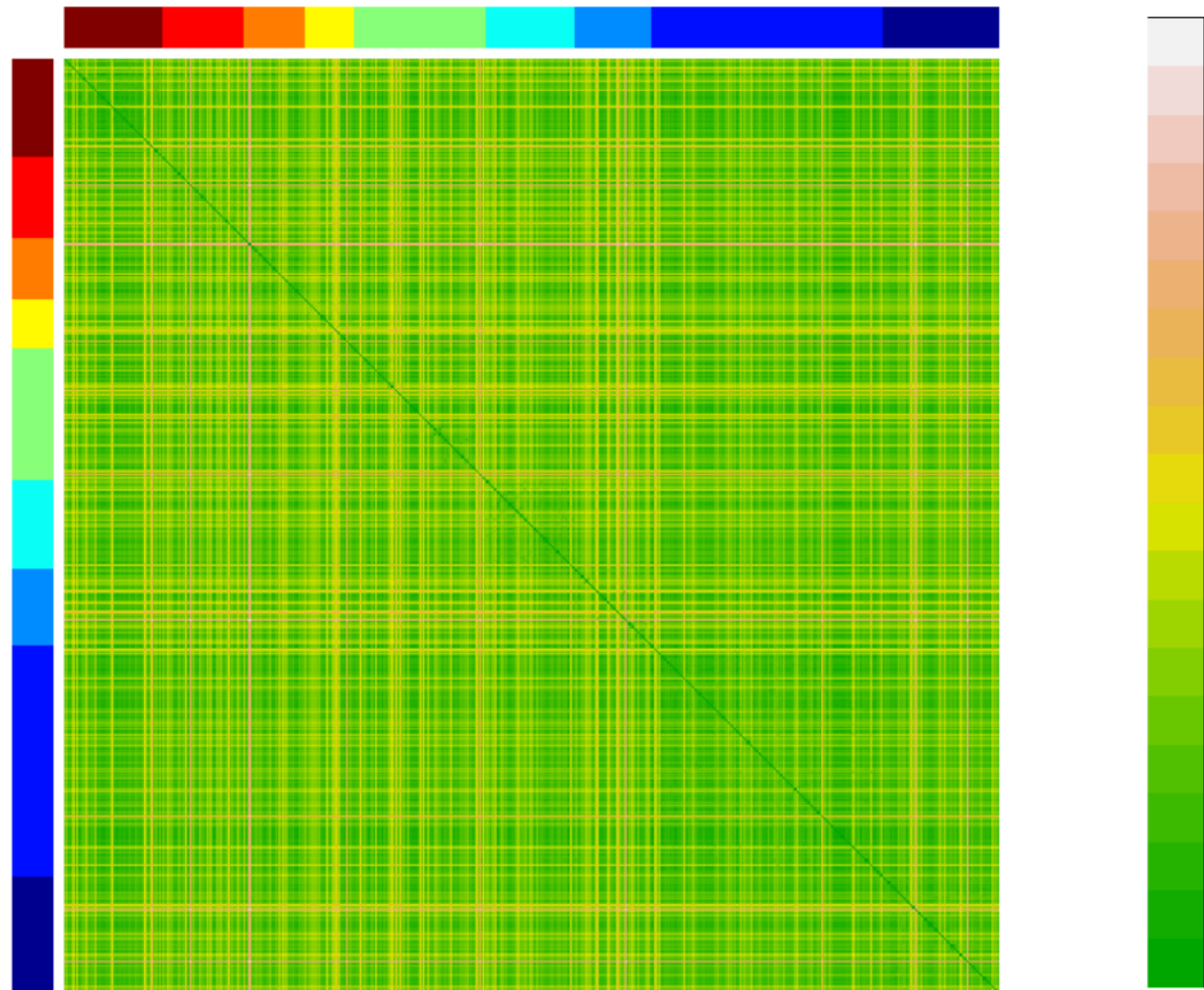
Minkowski, L=1
order by genera

Manhattan distance heatmap



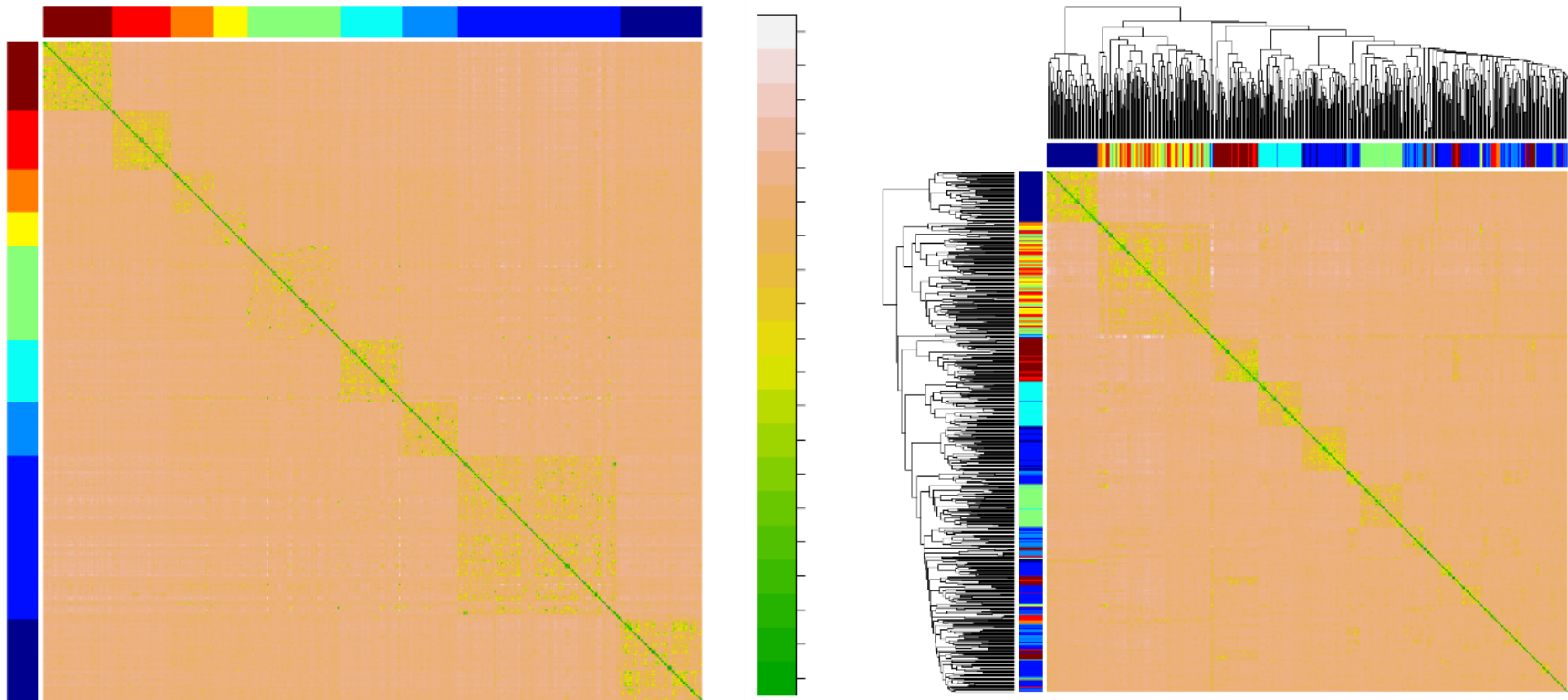
Minkowski, L=2
order by genera

Euclidean distance heatmap



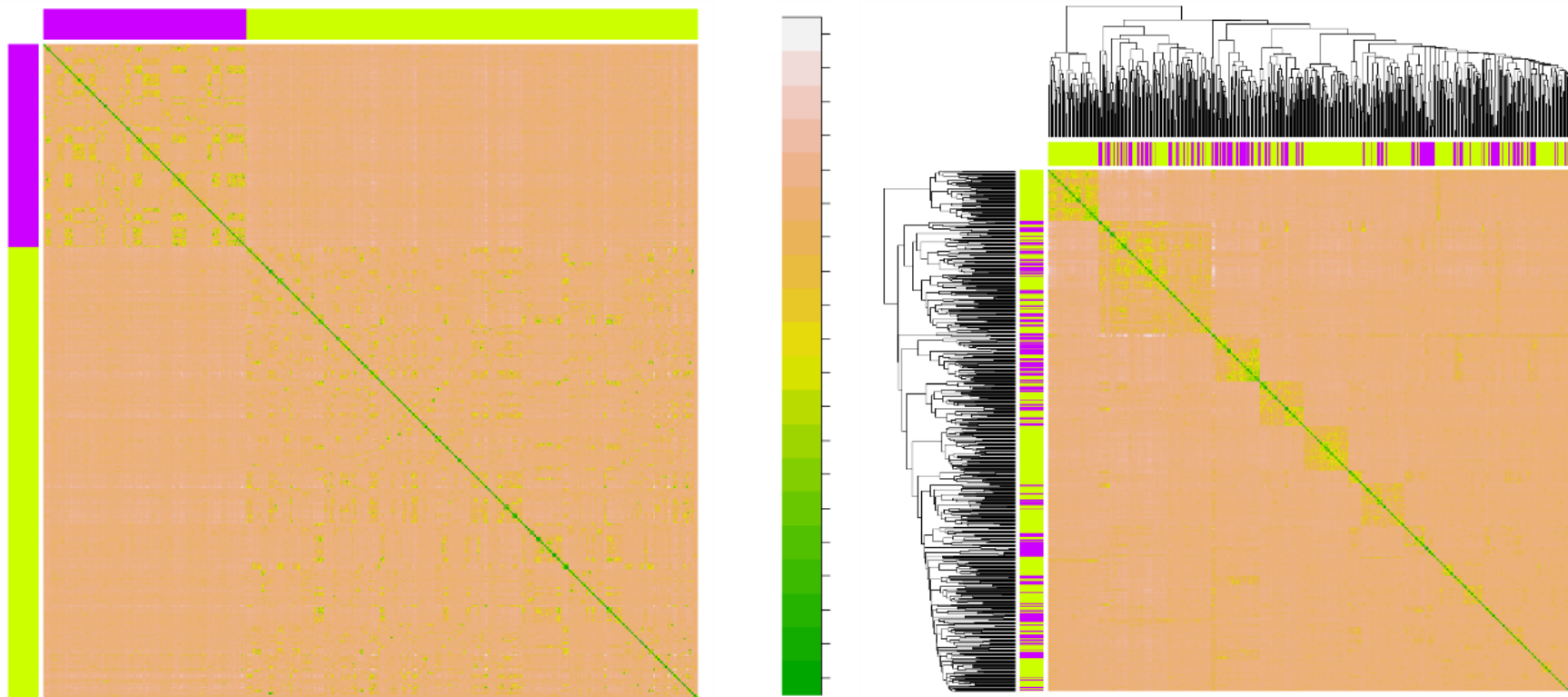
Correlation-based similarity heatmap

order by genera, clustering by average-link



Correlation-based similarity heatmap

order by gram type, clustering by average-link



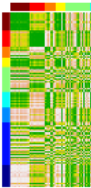
Dimension Reduction

ISOMap for 3-dimensional space

Shrinkage PCA for 3-dimensional space

MDS for 3-dimensional space

ISOMap with k=6
correlation-based
order by genera.

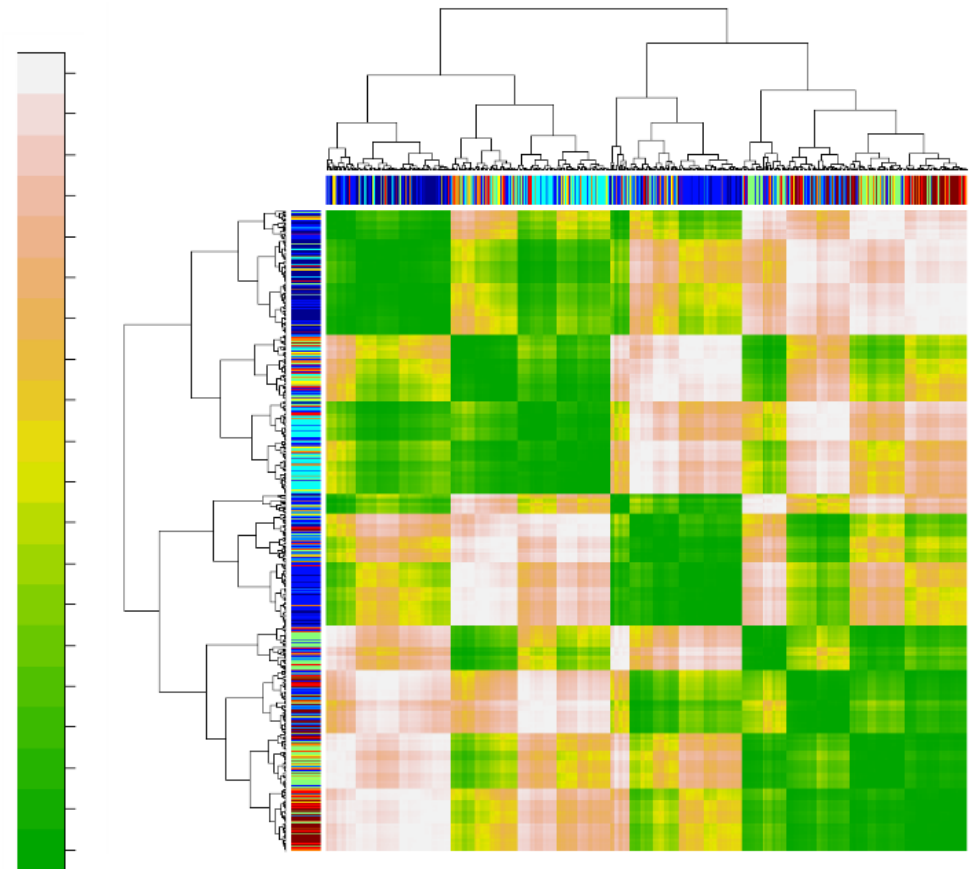
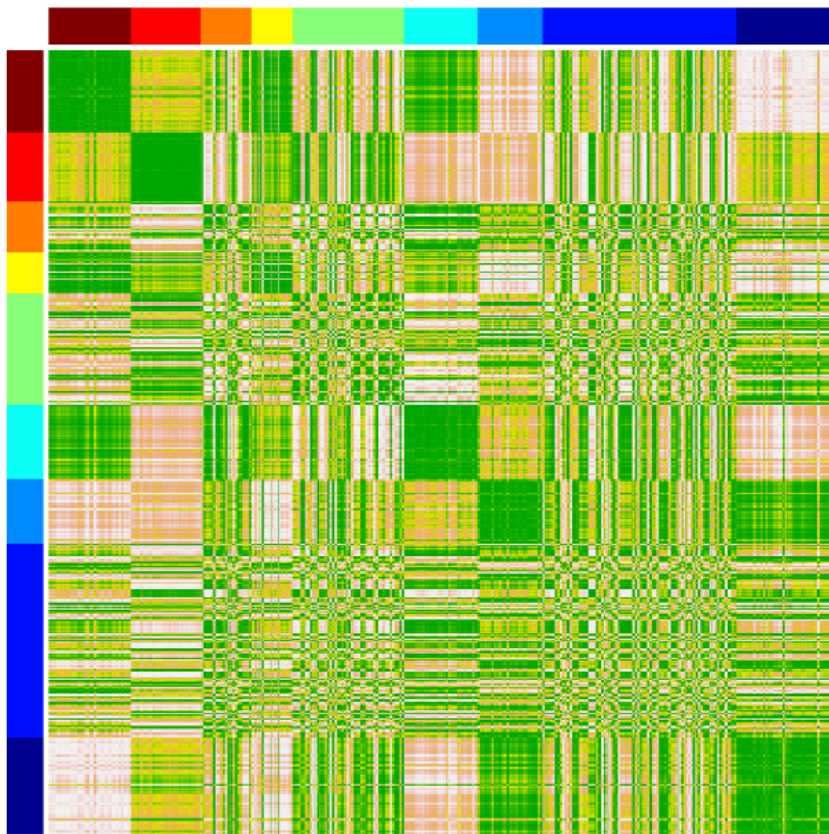


ISOMap

ISOMap with $k=5$

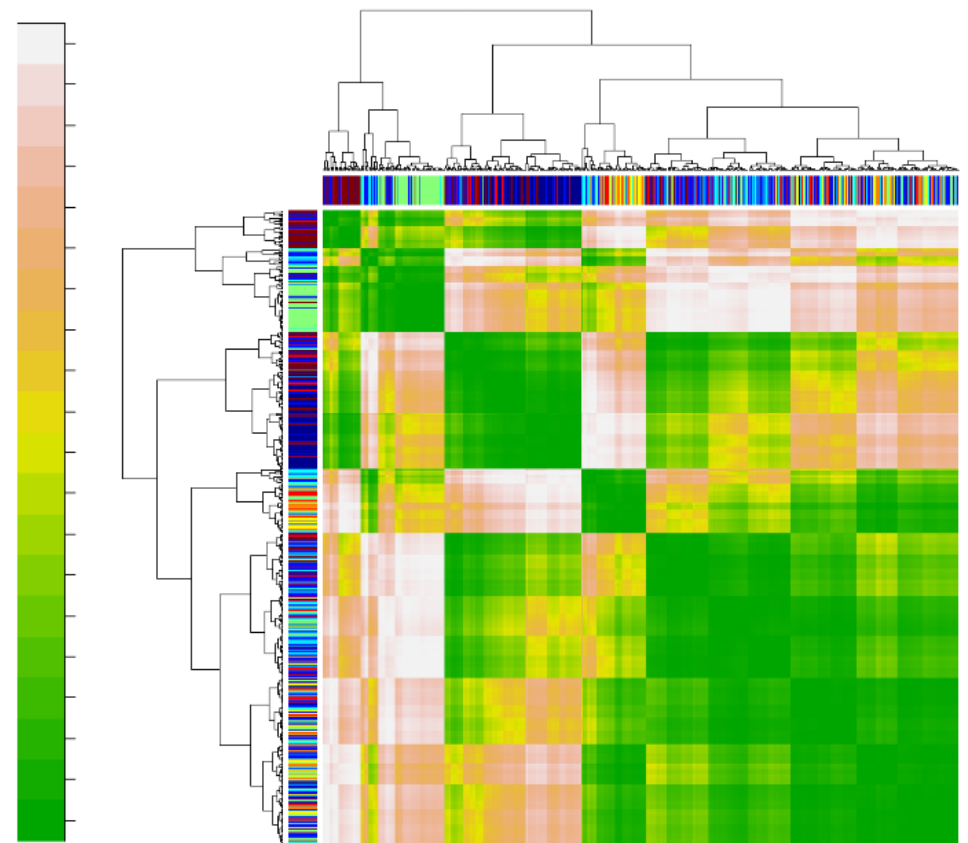
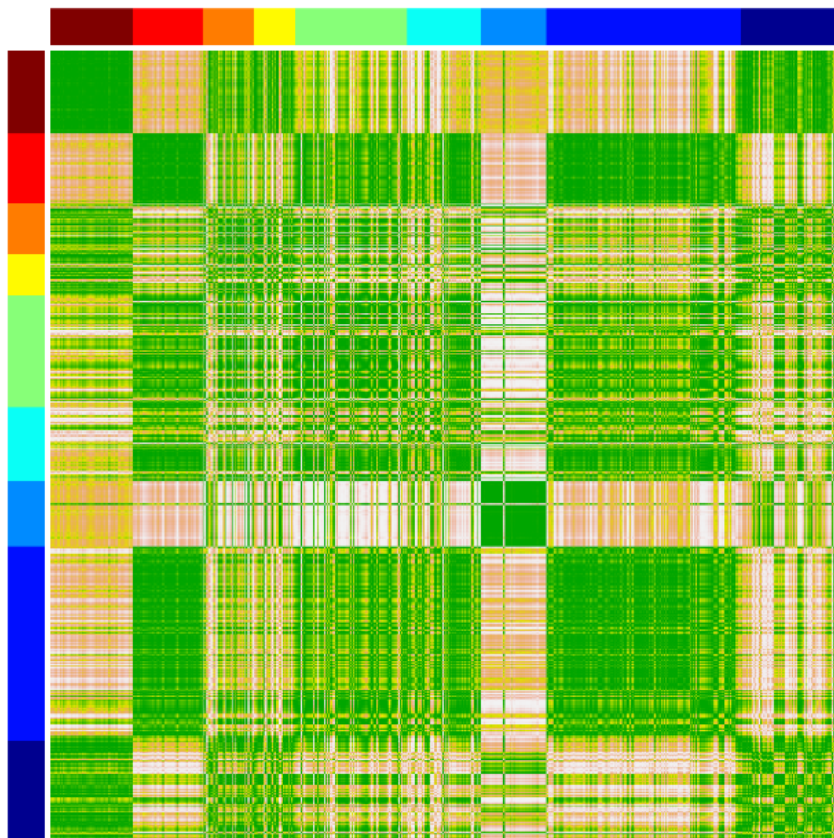
correlation-based similarity heatmap

order by genera, clustering by average-link



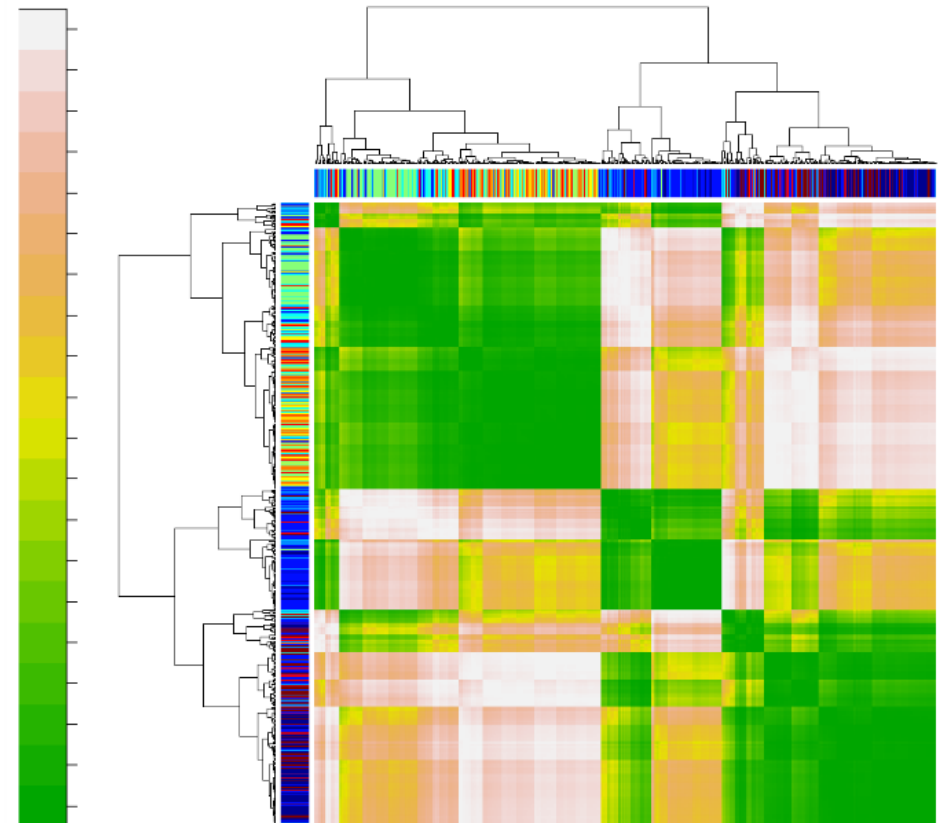
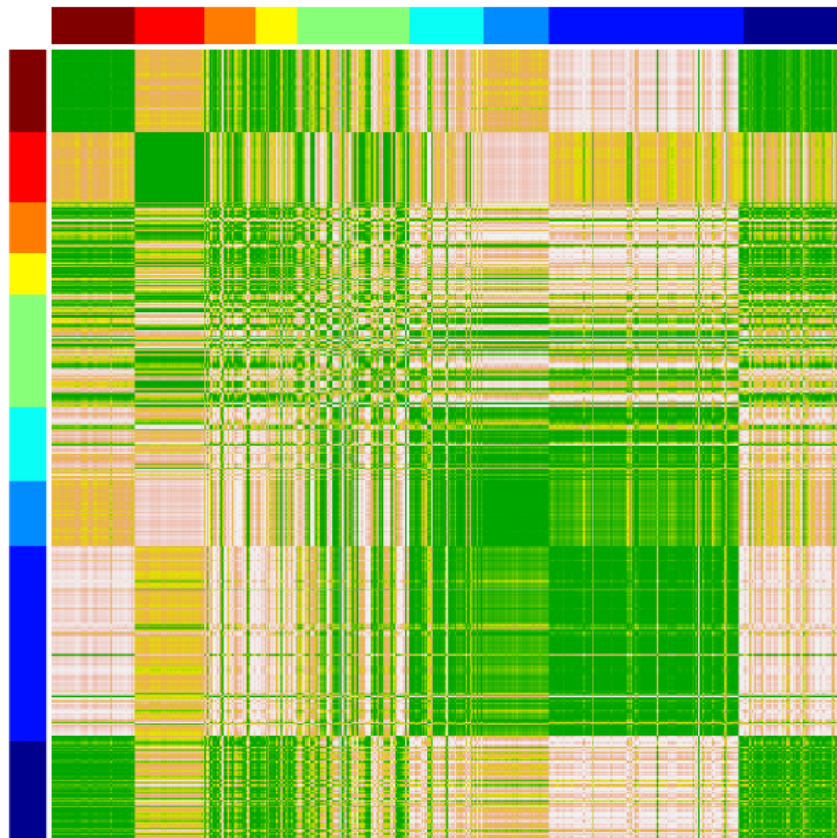
Shrinkage PCA

correlation-based similarity heatmap
order by genera, clustering by average-link



Multidimensional scaling

correlation-based similarity heatmap
order by genera, clustering by average-link





*Performance of
higher dimension selected*

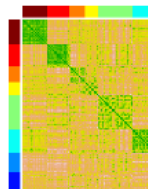
Shrinkage PCA for 30-dimension

ISOMap for 30-dimension

MDS for 30-dimension

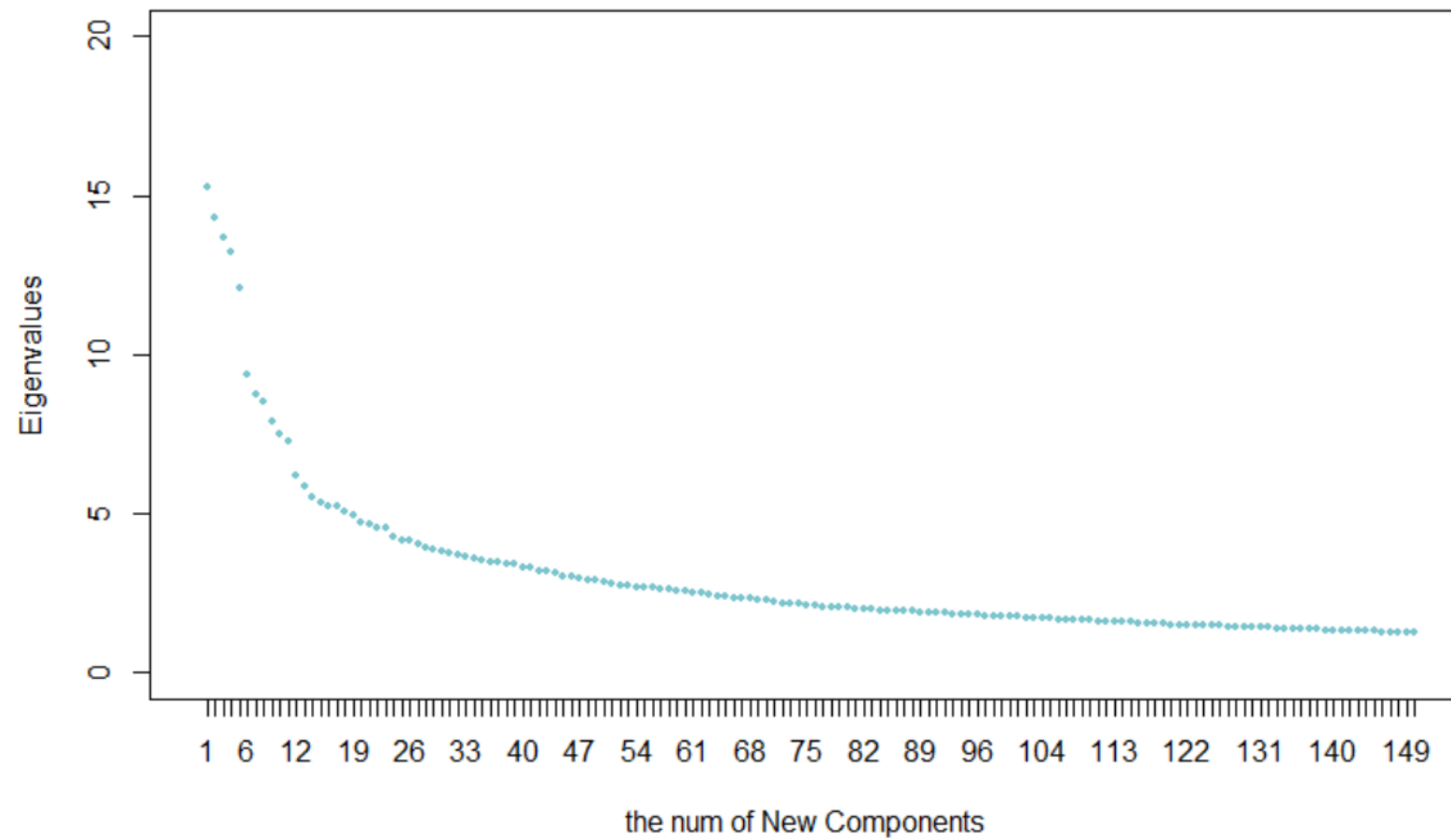


ISOMap with k=5
correlation-based
order by genera,



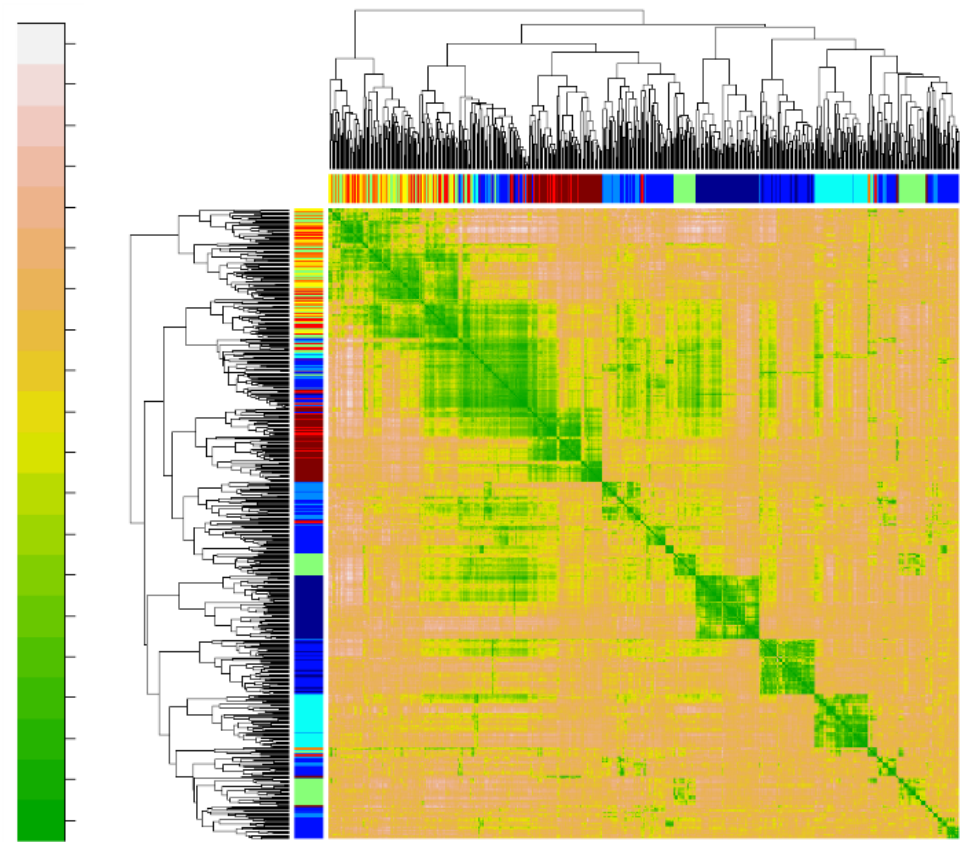
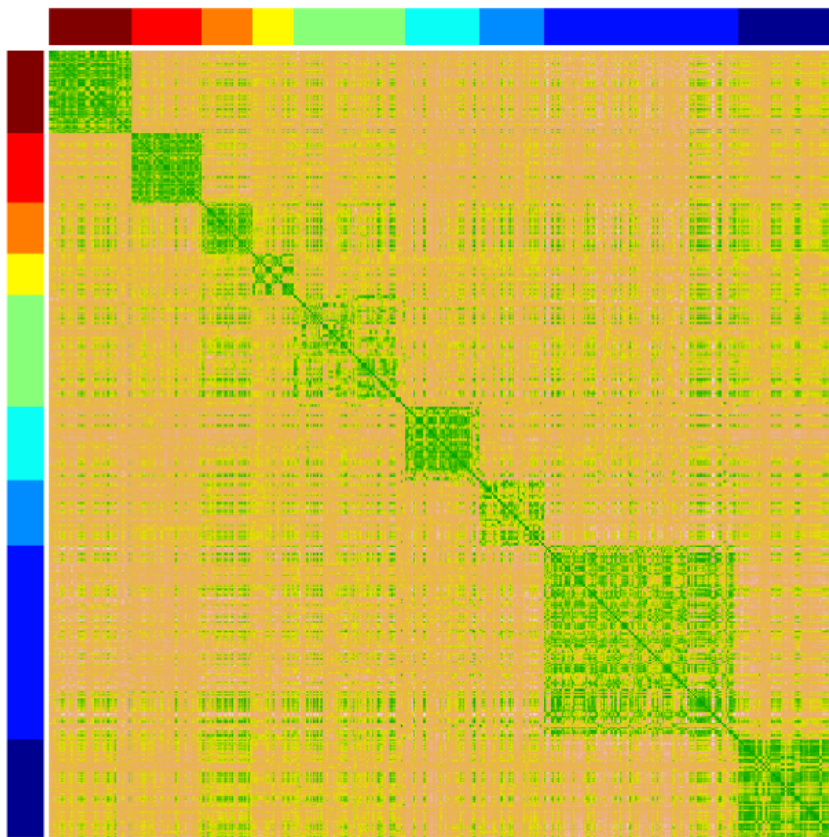
Shrinkage PCA

eigenvalues plot



Shrinkage PCA

correlation-based similarity heatmap
order by genera, clustering by average-link

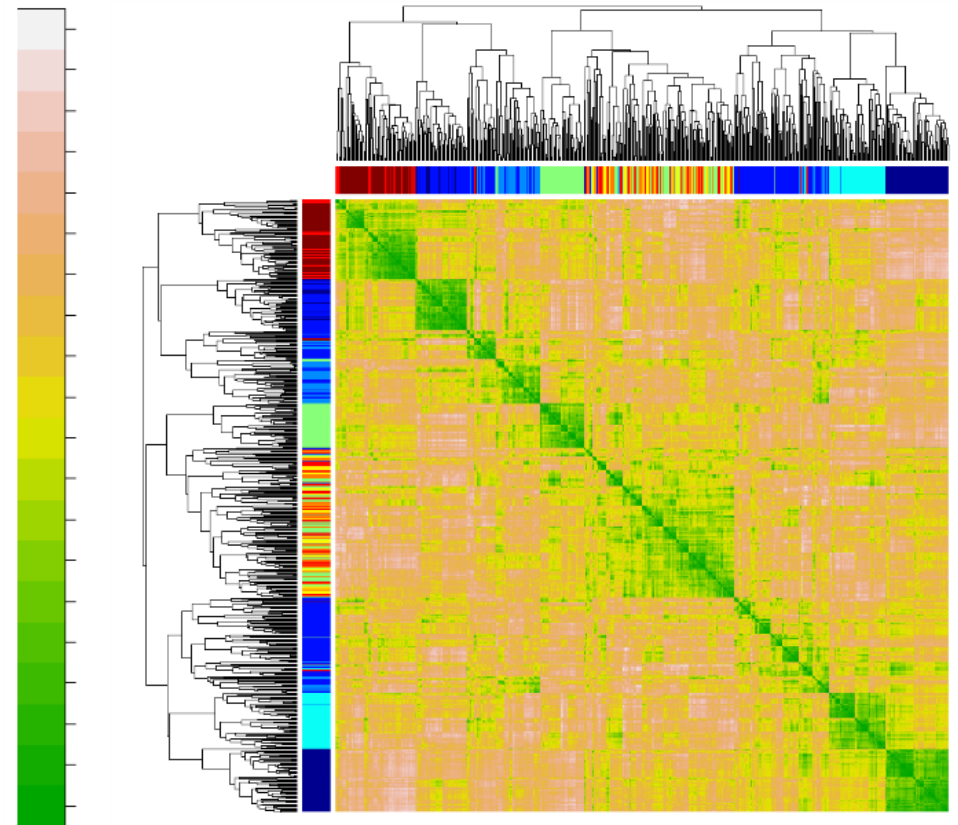
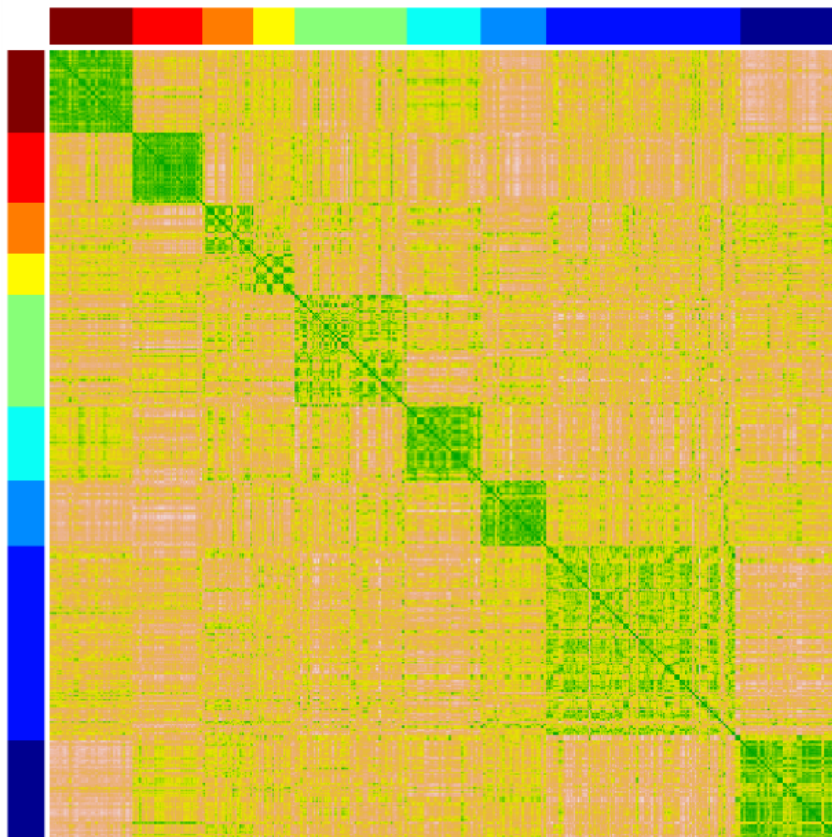


ISOmap

ISOmap with $k=5$

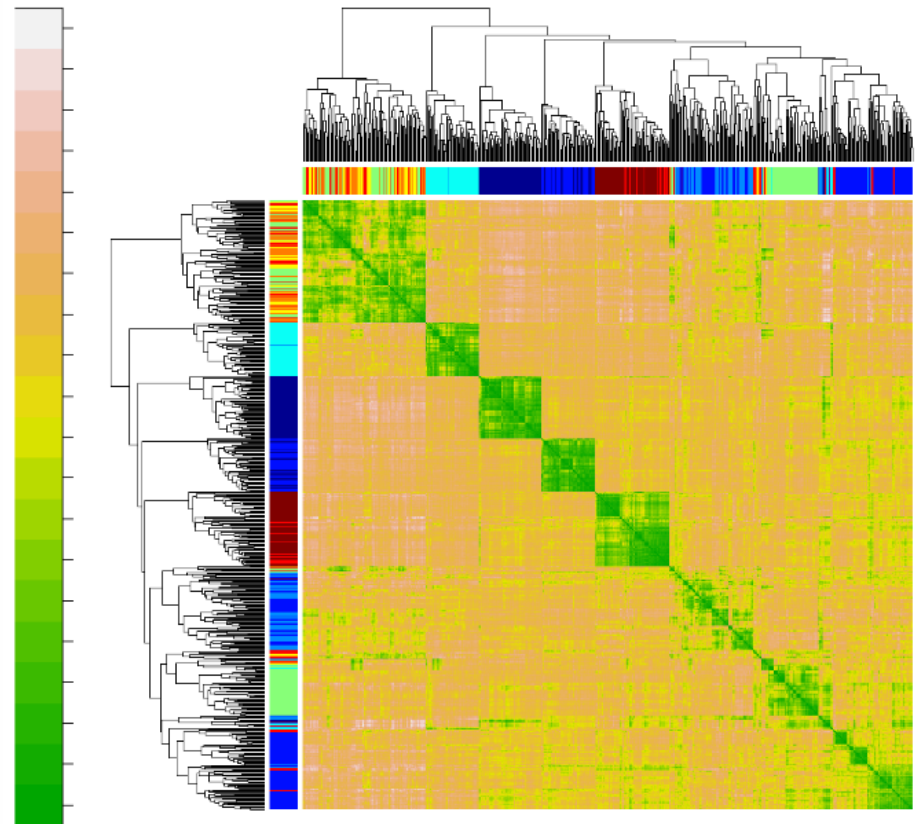
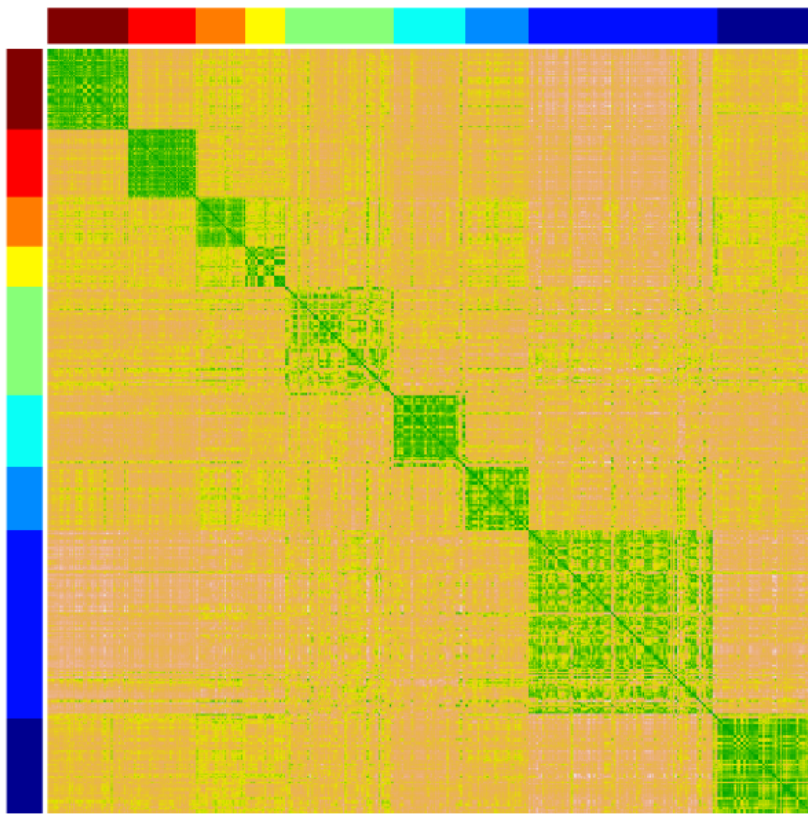
correlation-based similarity heatmap

order by genera, clustering by average-link



Multidimensional scaling

correlation-based similarity heatmap
order by genera, clustering by average-link





Classification

Support vector machine with kernel = **linear**

Subsampling the training dataset and the testing dataset with the ratio **2 : 1**

The accuracy of prediction is obtained from processing SVM for **50 times**

Stability of the classifying performance of SVM

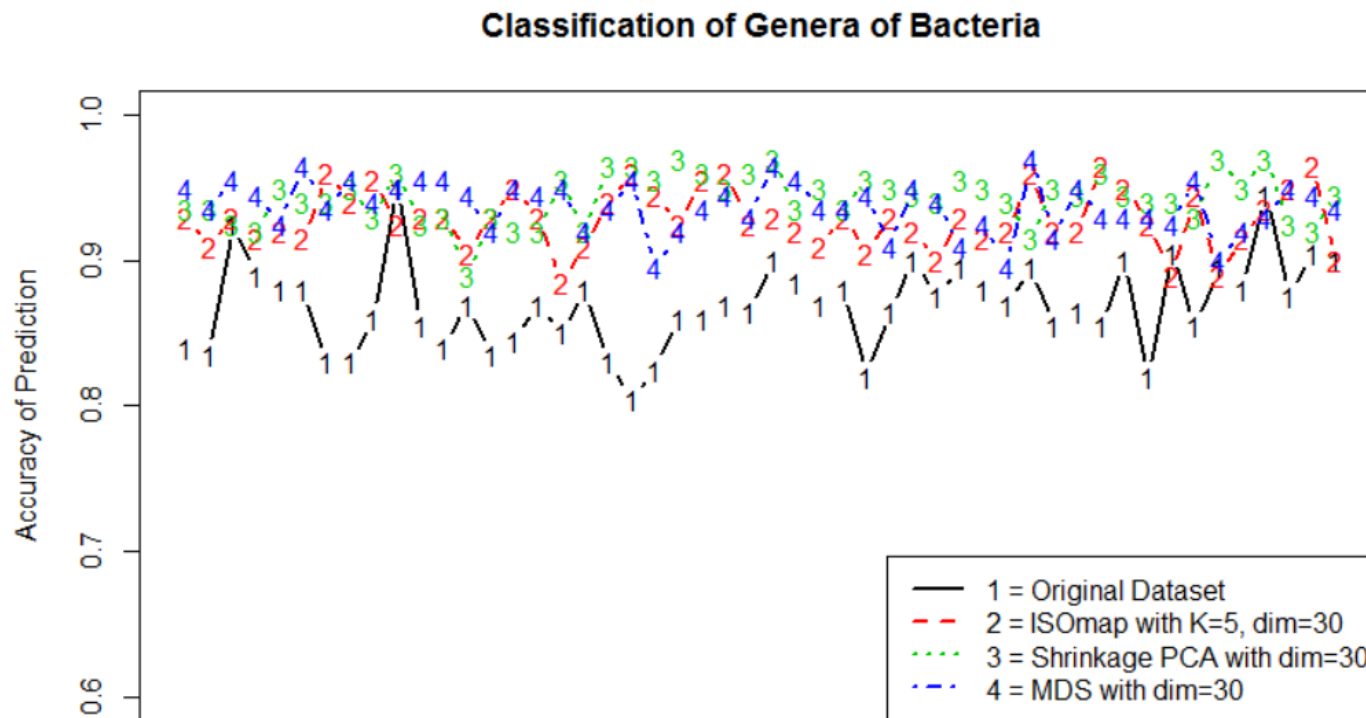
of Bacteria

Classification of Gram Type



Classification of Genera of Bacteria

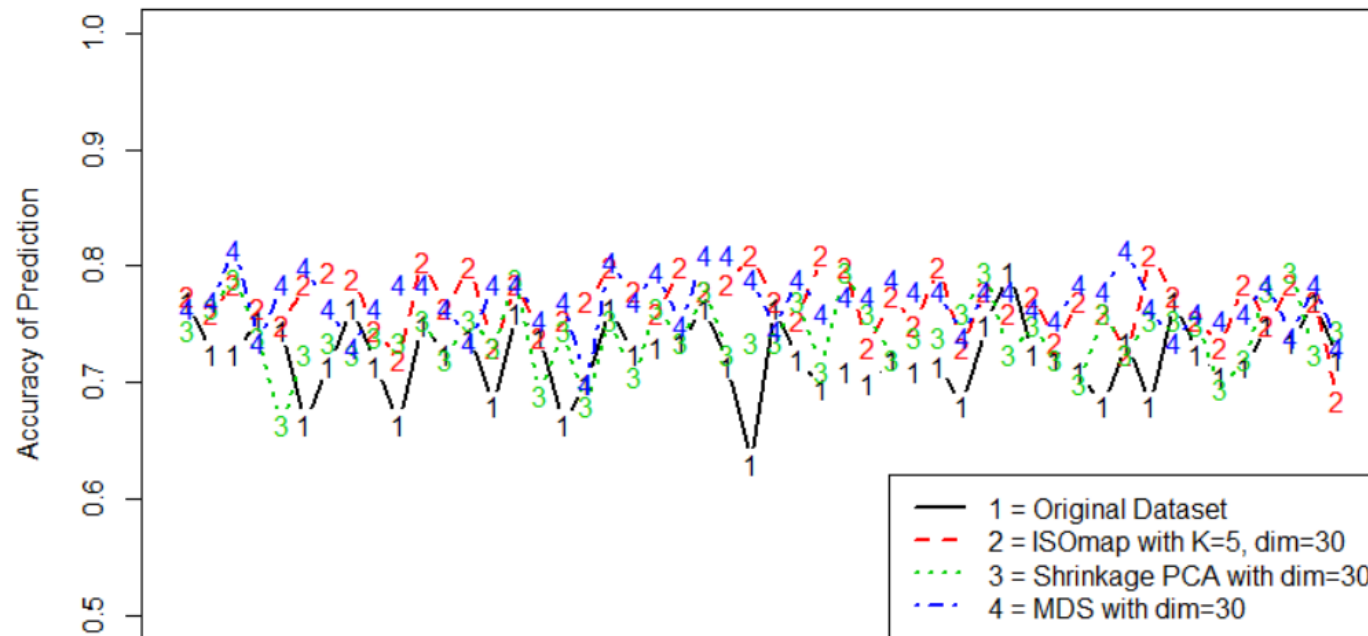
```
> colMeans(mat1)
[1] 0.8694 0.9284 0.9426 0.9372
> apply(mat1,2,function(x) round(var(x)^0.5,5))
[1] 0.03092 0.02026 0.01703 0.01762
```



Classification of Species of Bacteria

```
> colMeans(mat2)
[1] 0.7234 0.7681 0.7411 0.7707
> apply(mat2,2,function(x) round(var(x)^0.5,5))
[1] 0.03325 0.02723 0.02949 0.02443
```

Classification of Species of Bacteria



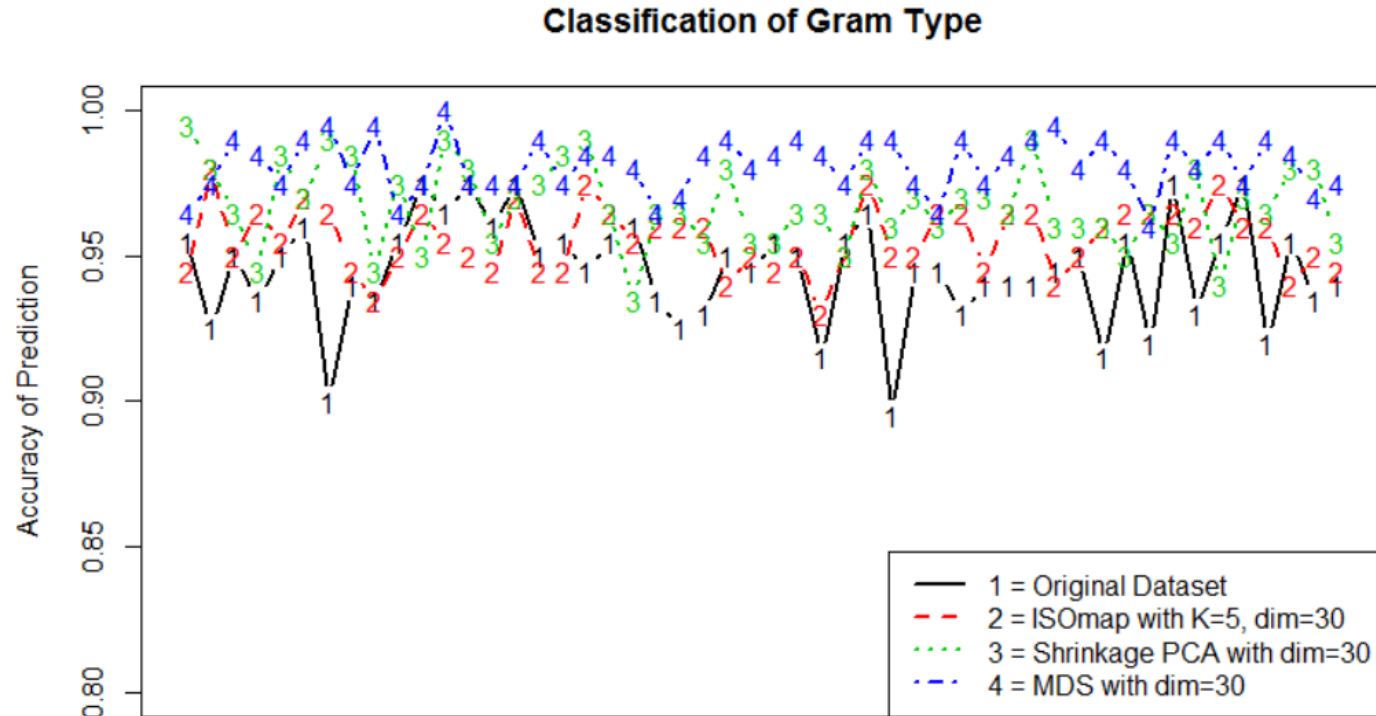
Classification of Gram Type

```
> colMeans(mat3)
```

```
[1] 0.9450 0.9558 0.9674 0.9814
```

```
> apply(mat3,2,function(x) round(var(x)^0.5,5))
```

```
[1] 0.01841 0.01131 0.01426 0.00937
```

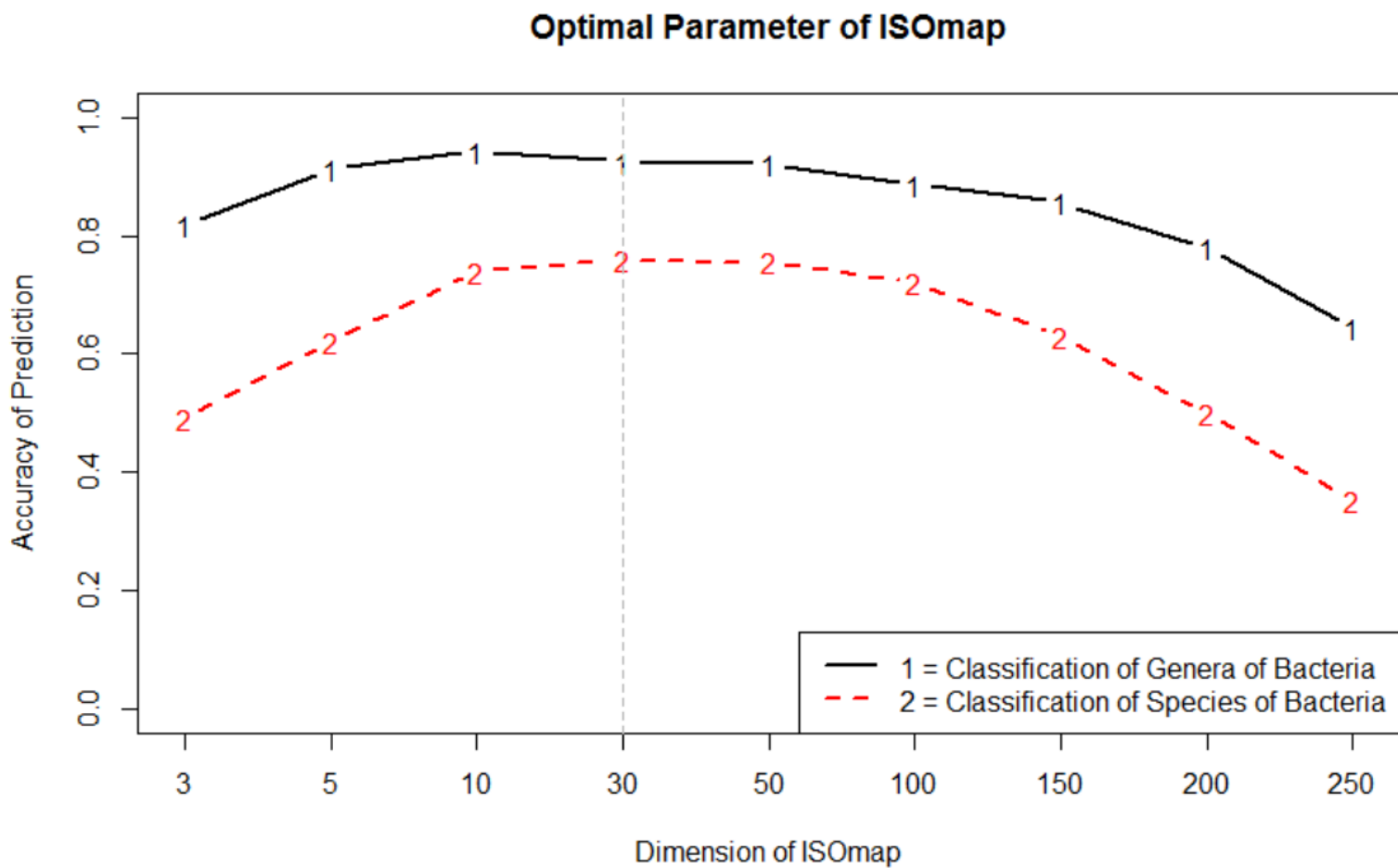


Discussing of Optimization

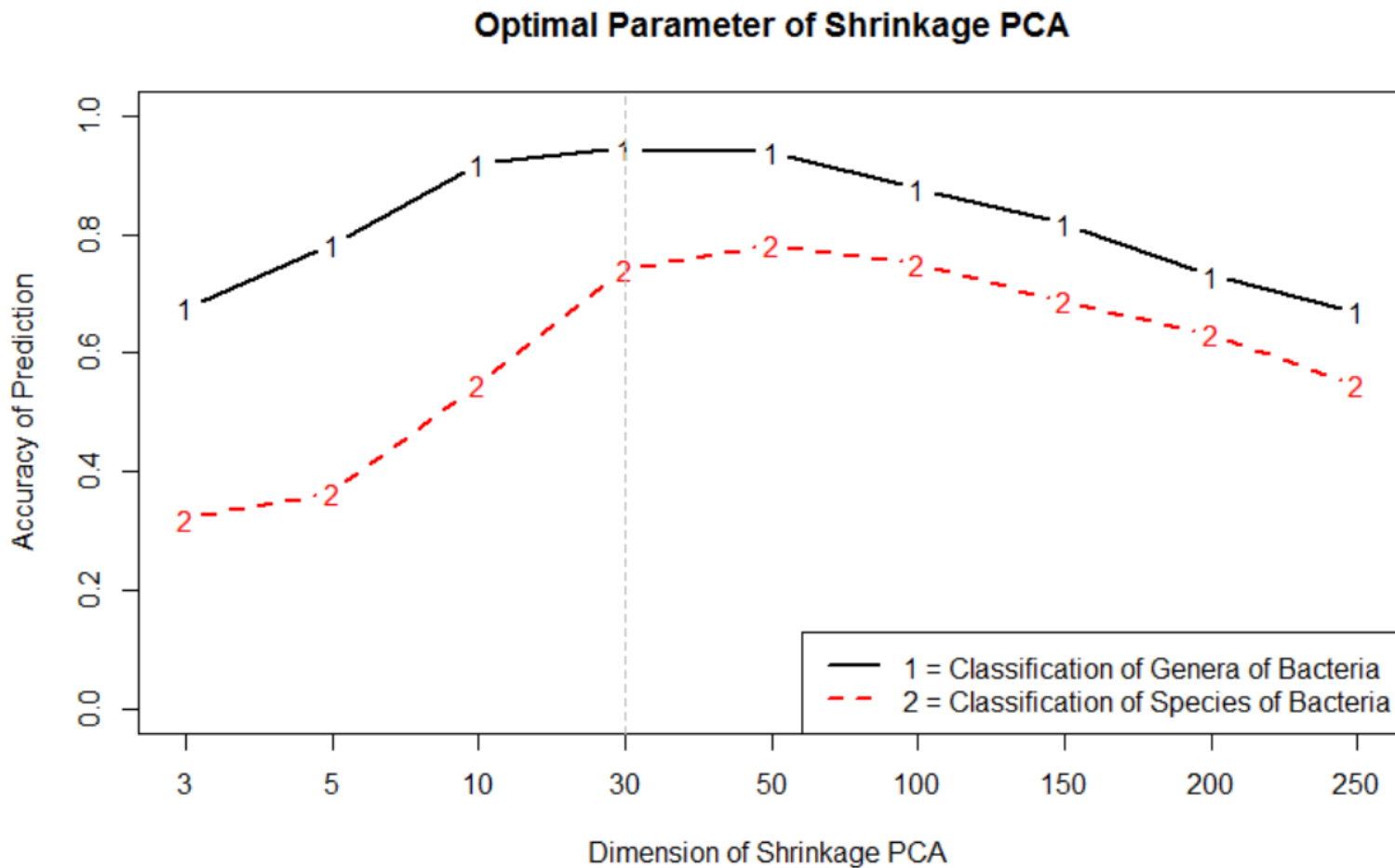
The performance about different numbers of dimension selected for models



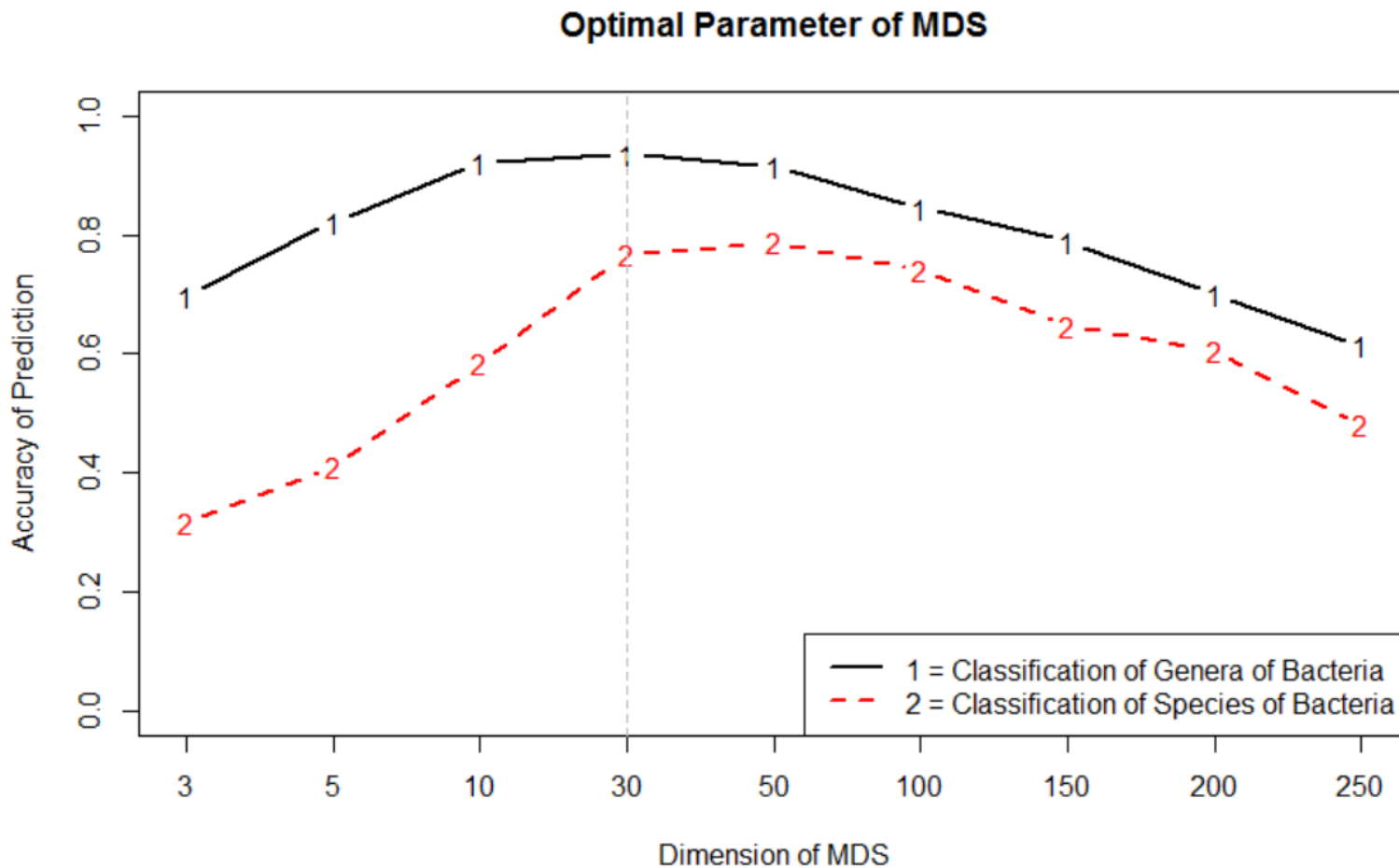
ISOMap



Shrinkage PCA



Multidimensional scaling



THE END

THANK YOU FOR LISTENING