

# Ionflow: network and enrichment analysis for ionomics data

*Wanchang Lin*

2020-12-21

## Contents

Data preparation . . . . .	2
Data pre-processing . . . . .	2
Data filtering . . . . .	7
Data clustering . . . . .	7
Gene network. . . . .	8
Enrichment analysis . . . . .	14
Exploratory analysis . . . . .	17

## Ionflow: network and enrichment analysis for ionomics data

This vignette explains how to perform ionomics data analysis including gene network and enrichment analysis by using a modification of the R package, [ionflow](#). The modification([ionflow\\_funcs](#)) was made by Wanchang Lin ([w.lin@imperial.ac.uk](mailto:w.lin@imperial.ac.uk)) and Jacopo Iacovacci([j.iacovacci@imperial.ac.uk](mailto:j.iacovacci@imperial.ac.uk)).

### Data preparation

To explore the process, we'll use the ionomics data set:

```
ion_data <- read.table("./test-data/iondata.tsv", header = T, sep = "\t")
dim(ion_data)
#> [1] 9999 16
```

Ten random data records are shown as:

```
sample_n(ion_data, 10)
```

**Table 1: Samples of raw data**

Knockout	Batch_ID	Ca	Cd	Co	Cu	Fe	K	Mg	Mn	Mo	Na	Ni	P	S	Zn
YBR268W	7	30.98	0.85	0.18	1.35	4.83	895.83	470.48	0.61	0.47	130.76	0.93	2915.70	490.14	16.20
YDR338C	10	9.54	0.85	0.20	1.37	3.75	2125.86	443.39	0.90	0.58	120.38	1.27	3026.55	310.02	15.59
YHR184W	19	35.32	0.82	0.11	1.16	3.30	1735.01	541.96	1.04	1.14	157.05	0.81	3717.33	374.17	12.27
YLR396C	80	12.85	1.19	0.17	1.15	5.13	1642.45	832.25	0.75	0.44	112.16	1.06	5108.83	714.23	13.44
YKL054C	23	47.64	0.56	0.19	1.96	17.31	2746.18	606.05	0.88	0.61	91.20	1.20	4500.47	459.67	17.22
YGL237C	87	45.11	1.08	0.15	1.56	8.83	2163.45	689.67	0.82	0.67	374.06	1.03	4222.77	555.38	12.78
YOR200W	80	9.71	0.89	0.16	1.09	5.38	685.06	521.40	0.61	1.89	219.95	0.62	3802.76	502.76	13.66
YDL227C	90	36.52	0.79	0.13	1.75	6.49	2667.00	731.12	1.31	0.89	244.10	0.85	4565.24	571.46	13.36
YHR161C	19	39.97	0.87	0.12	1.35	5.26	1876.79	536.72	0.94	1.17	131.28	0.89	3771.25	362.43	14.15
YGR133W	14	60.88	1.08	0.18	2.15	7.52	3966.02	802.69	1.59	0.98	233.48	1.59	5303.77	470.84	22.05

The first few columns are meta information such as gene ORF and batch id. The rest is the ionomics data.

### Data pre-processing

The raw data set should be pre-processed. The pre-processing function `PreProcessing` has functions:

- log transformation
- batch correction
- outlier detection
- standardisation

The raw data are at first log transformed and then followed by the batch correction. User can chose not to perform batch correction, otherwise default will be either *median* or *median* plus *std* method. If there is quality control for the batch correction, the user can use it and indicates in the argument of `control_lines`. Also one argument gives

## Ionflow: network and enrichment analysis for ionomics data

the user the option on how to use these control lines (`control_use`): If `control_use` is `control`, these control lines (data rows) are used for the batch correction factor; if `control.out`, others lines are used.

This data set has a control line: **YDL227C** mutant. The code segment below is to identify it:

```
max(with(ion_data, table(Knockout)))
#> [1] 1617
which.max(with(ion_data, table(Knockout)))
#> YDL227C
#>      209
```

The next stage is outlier detection. Here only univariate methods are implemented, including *mad*, *IQR*, and *log.FC.dist*. And like batch correction, the user can skip this procedure by setting `method_outliers = none` in the function argument. There is a threshold to control the number of outliers. The larger the threshold (`thres_outl`) the more outlier removal.

Standardisation provides three methods: *std*, *mad* or *custom*. If the method is *custom*, the user uses a specific *std* file like:

```
std <- read.table("./test-data/user_std.tsv", header = T, sep = "\t")
std
#>      Ion      sd
#> 1  Ca 0.1508
#> 2  Cd 0.0573
#> 3  Co 0.0580
#> 4  Cu 0.0735
#> 5  Fe 0.1639
#> 6   K 0.0940
#> 7  Mg 0.0597
#> 8  Mn 0.0771
#> 9  Mo 0.1142
#> 10 Na 0.1075
#> 11 Ni 0.0784
#> 12  P 0.0597
#> 13  S 0.0801
#> 14 Zn 0.0671
```

The pre-processing procedure returns not only processed ionomics data but also a symbolic data set. This data set is based on the ionomics data and is determined by a `threshold(thres_symb)`:

- 0 if ionomics value is located in `[-thres_symb, thres_symb]`
- 1 if ionomics value is larger than `thres_symb`
- -1 if ionomics value is smaller than `-thres_symb`

## Ionflow: network and enrichment analysis for ionomics data

Note that the symbolic data set is important since the key part of the network and enrichment analysis is based on the hierarchical clustering of symbolic data.

Let's run the pre-process procedure:

```
pre <- PreProcessing(data = ion_data,
  var_id = 1, batch_id = 2, data_id = 3,
  method_norm = "median",
  control_lines = "YDL227C",
  control_use = "control",
  method_outliers = "IQR",
  thres_outl = 3,
  stand_method = "std",
  stdev = NULL,
  thres_symb = 3)

names(pre)
#> [1] "stats.raw_data"      "stats.outliers"      "stats.batch_data"
#> [4] "data.long"           "data.gene.logFC"     "data.gene.zscores"
#> [7] "data.gene.symb"      "plot.dot"            "plot.hist"
```

The results include summaries of raw data and processed data. The latter is:

```
pre$stats.batch_data %>%
  kable(caption = 'Processed data summary', digits = 2, booktabs = T) %>%
  kable_styling(full_width = F, font_size = 10)
```

**Table 2: Processed data summary**

Ion	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Variance
Ca	-4.45	-0.28	-0.13	-0.12	0.02	2.35	0.11
Cd	-1.70	0.03	0.10	0.11	0.17	0.93	0.03
Co	-2.80	0.02	0.09	0.06	0.15	1.60	0.05
Cu	-0.66	-0.10	-0.03	-0.01	0.04	5.28	0.04
Fe	-7.48	-0.17	-0.06	-0.02	0.07	6.88	0.14
K	-2.21	-0.17	-0.01	-0.08	0.09	1.83	0.08
Mg	-1.84	-0.06	0.01	-0.01	0.07	1.69	0.03
Mn	-4.11	-0.24	-0.08	-0.13	0.01	1.78	0.06
Mo	-2.03	-0.26	-0.08	-0.08	0.09	4.44	0.13
Na	-7.41	-0.53	-0.22	-0.33	-0.04	1.25	0.24
Ni	-2.40	-0.01	0.09	0.12	0.21	7.90	0.12
P	-1.18	-0.06	0.00	-0.01	0.06	1.45	0.02
S	-2.38	-0.03	0.05	0.06	0.16	2.38	0.04
Zn	-0.46	-0.08	-0.03	-0.01	0.03	4.60	0.02

The pre-processed data and symbolic data are like this:

## Ionflow: network and enrichment analysis for ionomics data

```
pre$data.gene.zscores %>% head() %>%  
  kable(caption = 'Processed data', digits = 2, booktabs = T) %>%  
  kable_styling(full_width = F, font_size = 10,  
                latex_options = c("striped", "scale_down"))
```

**Table 3: Processed data**

Line	Ca	Cd	Co	Cu	Fe	K	Mg	Mn	Mo	Na	Ni	P	S	Zn
YAL004W	-1.16	0.75	1.19	-0.47	0.04	0.61	0.51	-0.84	-0.08	-1.84	1.71	0.52	0.33	-0.09
YAL005C	-1.67	0.84	0.55	0.58	-2.79	0.59	0.31	-1.16	-1.42	-0.12	1.48	0.73	0.13	-0.13
YAL007C	-2.12	0.64	0.23	-0.53	-0.24	0.79	-0.09	-0.14	1.22	-0.92	0.00	0.09	-0.29	-0.65
YAL008W	-2.34	1.13	0.21	-0.73	-2.16	0.52	-0.02	-0.87	0.93	-0.58	0.02	-0.09	-0.73	-0.47
YAL009W	-1.18	0.66	0.55	-1.11	-3.91	0.22	0.09	-0.18	1.50	-0.84	-0.09	0.14	0.01	-0.36
YAL010C	-1.28	1.43	2.27	0.46	1.53	-2.75	0.04	-0.74	-9.71	-4.30	2.42	-0.98	-0.05	-0.01

```
pre$data.gene.symb %>% head() %>%  
  kable(caption = 'Symbolic data', booktabs = T) %>%  
  kable_styling(full_width = F, font_size = 10)
```

**Table 4: Symbolic data**

Line	Ca	Cd	Co	Cu	Fe	K	Mg	Mn	Mo	Na	Ni	P	S	Zn
YAL004W	0	0	0	0	0	0	0	0	0	0	0	0	0	0
YAL005C	0	0	0	0	0	0	0	0	0	0	0	0	0	0
YAL007C	0	0	0	0	0	0	0	0	0	0	0	0	0	0
YAL008W	0	0	0	0	0	0	0	0	0	0	0	0	0	0
YAL009W	0	0	0	0	-1	0	0	0	0	0	0	0	0	0
YAL010C	0	0	0	0	0	0	0	0	-1	-1	0	0	0	0

The symbolic data are calculated from the processed data with control of `thres_symb` (here it is 3). You can obtain a new symbol data set by re-assigning a new threshold to the function `symbol_data`:

```
data_symb <- symbol_data(pre$data.gene.zscores, thres_symb = 2)  
data_symb %>% head() %>%  
  kable(caption = 'Symbolic data with threshold of 2', booktabs = T) %>%  
  kable_styling(full_width = F, font_size = 10)
```

The `thres_symb` is a crucial value to get the symbolic data. Before re-setting this threshold, the user should check the summary of processed data and pay attention to the maximum values. For example, some ions (for example, *Cd* and *Mn*) are all zero even with 2 of `thres_symb`.

The pre-processed data distribution is:

Ionflow: network and enrichment analysis for ionomics data

Table 5: Symbolic data with threshold of 2

Line	Ca	Cd	Co	Cu	Fe	K	Mg	Mn	Mo	Na	Ni	P	S	Zn
YAL004W	0	0	0	0	0	0	0	0	0	0	0	0	0	0
YAL005C	0	0	0	0	-1	0	0	0	0	0	0	0	0	0
YAL007C	-1	0	0	0	0	0	0	0	0	0	0	0	0	0
YAL008W	-1	0	0	0	-1	0	0	0	0	0	0	0	0	0
YAL009W	0	0	0	0	-1	0	0	0	0	0	0	0	0	0
YAL010C	0	0	1	0	0	-1	0	0	-1	-1	1	0	0	0

pre\$plot.hist

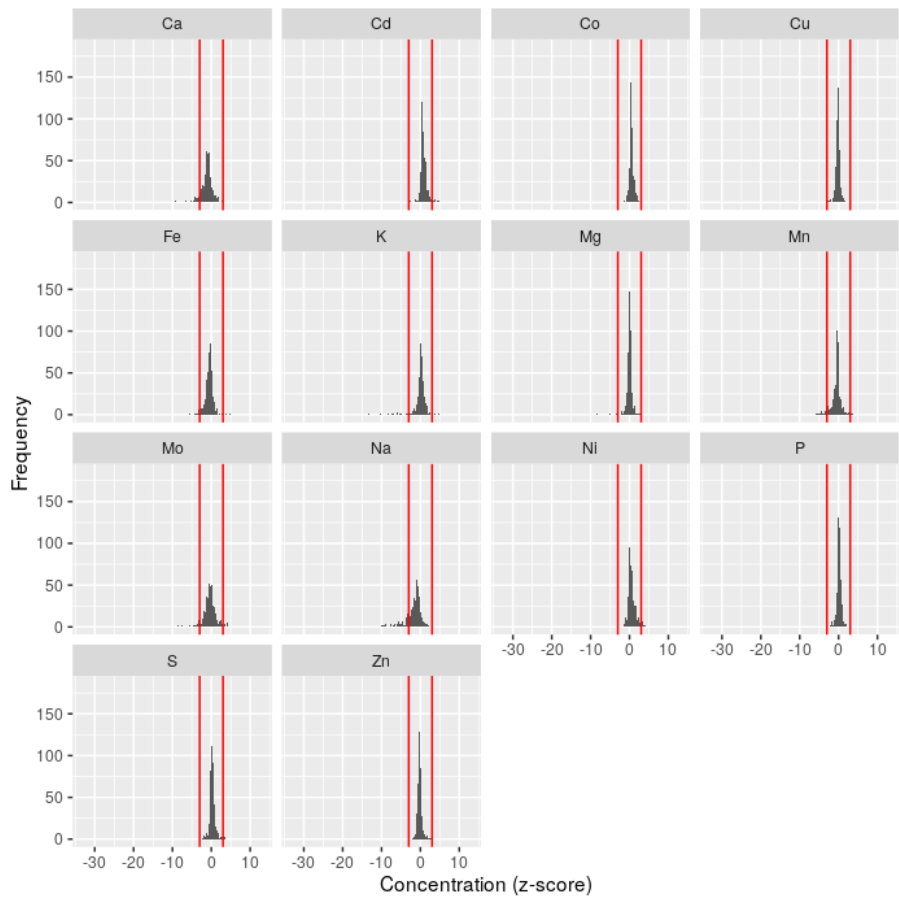


Figure 1: Ionomics data distribution plot

## Ionflow: network and enrichment analysis for ionomics data

### Data filtering

There are a lot of ways to filter genes. Here genes are filtered according to symbolic data: remove genes with all values which are zero.

```
data <- pre$data.gene.zscores
data_symb <- pre$data.gene.symb
idx <- rowSums(abs(data_symb[, -1])) > 0
dat <- data[idx, ]
dat_symb <- data_symb[idx, ]
dim(dat)
#> [1] 549 15
```

### Data clustering

The hierarchical cluster analysis is the key part of gene network and gene enrichment analysis. The methodology is as follow:

- Compute the distance of symbolic data
- Hierarchical cluster analysis on the distance
- Identify clusters/groups with a threshold of minimal number of cluster size

One example is:

```
clust <- gene_clus(dat_symb[, -1], min_clust_size = 10)
names(clust)
#> [1] "clus"      "idx"      "tab"      "tab_sub"
```

The cluster centres are:

```
clust$tab_sub
#>   cluster nGenes
#> 1      4     149
#> 2     11      72
#> 3      7      36
#> 4      1      27
#> 5     18      15
#> 6      5      12
#> 7      3      11
#> 8      8      11
```

This shows clusters and the number of genes (larger than `min_cluster_size`).

The identified gene located in those clusters are:

```
sum(clust$idx)                                #' numbers of all genes
#> [1] 333
```

## Ionflow: network and enrichment analysis for ionomics data

```
head(as.character(dat[,1][clust$idx])) #' and some are  
#> [1] "YAL009W" "YAL013W" "YAL014C" "YAL020C" "YAL021C" "YAL022C"
```

### Gene network

The gene network uses both the ionomics and symbolic data. The similarity measures on ionomics data are used to construct the network. Before creating a network, these analyses are further filtered by:

- clustering of symbolic data;
- and the similarity threshold located between 0 and 1;

The methods implemented are: *pearson*, *spearman*, *kendall*, *cosine*, *mahal\_cosine* or *hybrid\_mahal\_cosine*. The first three methods are correlation methods and *cosine* is similar to the Pearson correlation which is the [cosine similarity between two centred vectors](#). For the last two methods, see publication: [Extraction and Integration of Genetic Networks from Short-Profile Omic Data Sets](#) for details.

For example, we use the Pearson correlation as similarity measure for network analysis:

```
net <- GeneNetwork(data = dat,  
                   data_symb = dat_symb,  
                   min_clust_size = 10,  
                   thres_corr = 0.75,  
                   method_corr = "pearson")
```

The network with nodes coloured by the symbolic data clustering is:

```
net$plot.pnet1
```

The same network, but nodes are coloured by the network community detection:

```
net$plot.pnet2
```

The network analysis also returns a network impact and betweenness plot:

```
net$plot.impact_betweenness
```



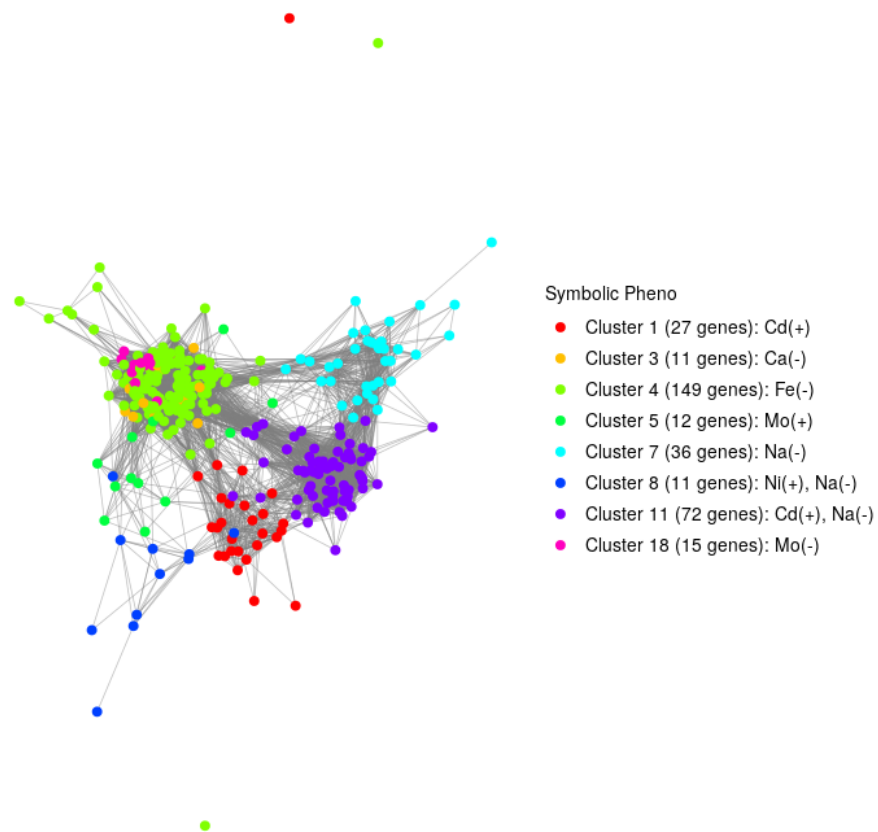


Figure 2: Network with Pearson correlation: symbolic clustering

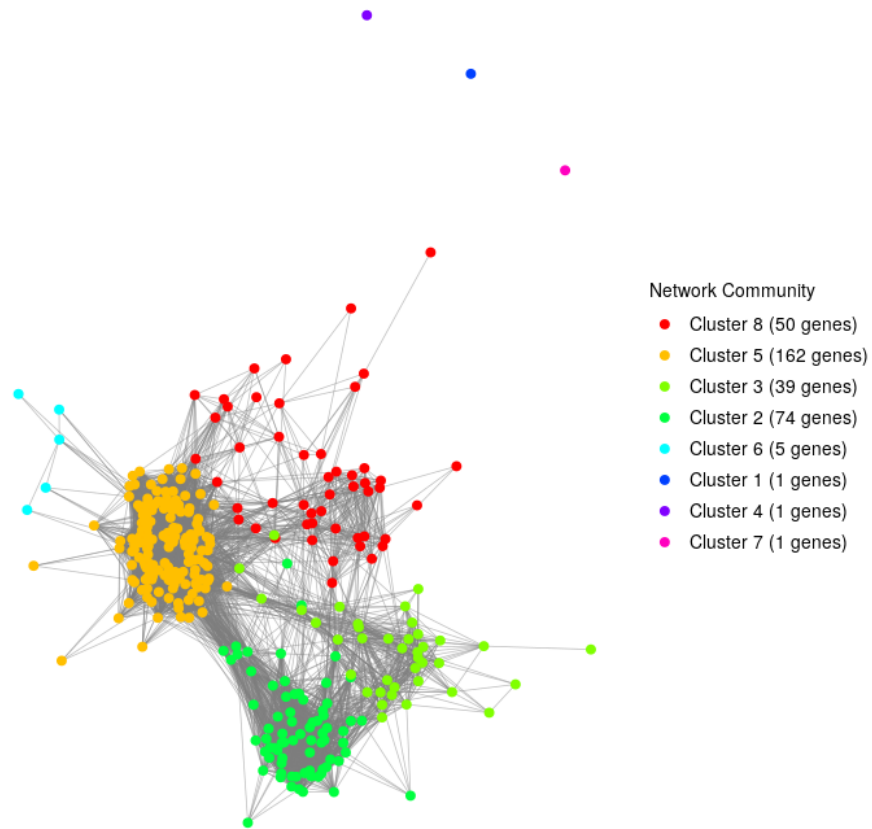


Figure 3: Network with Pearson correlation: community detection

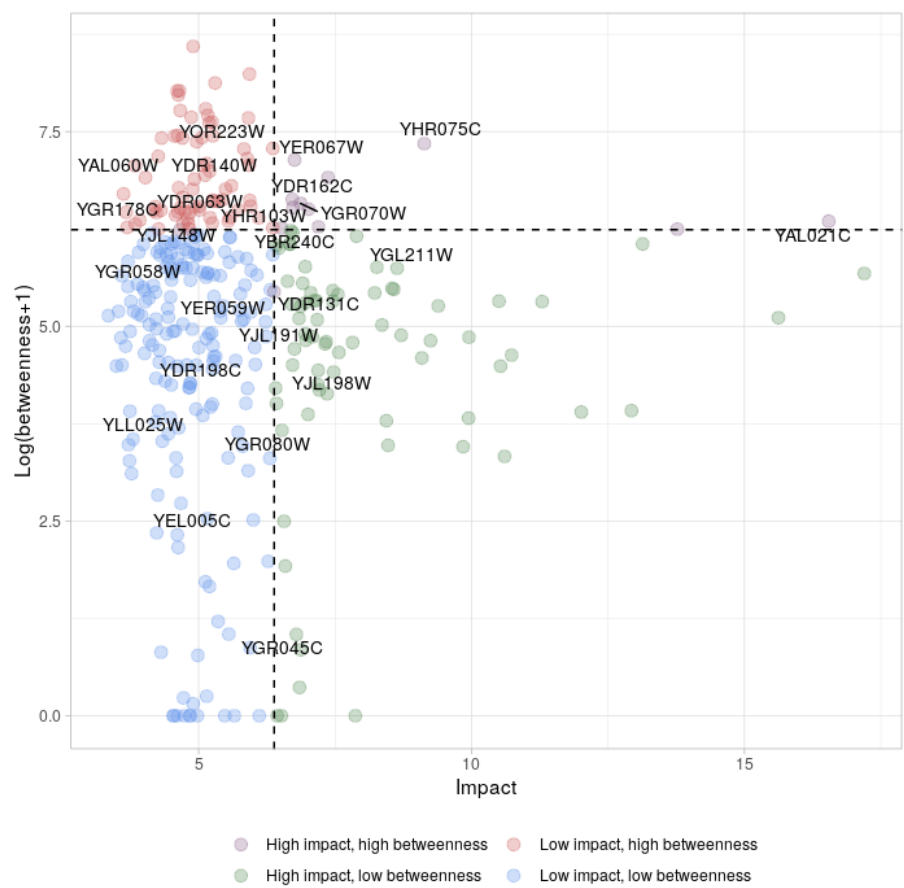
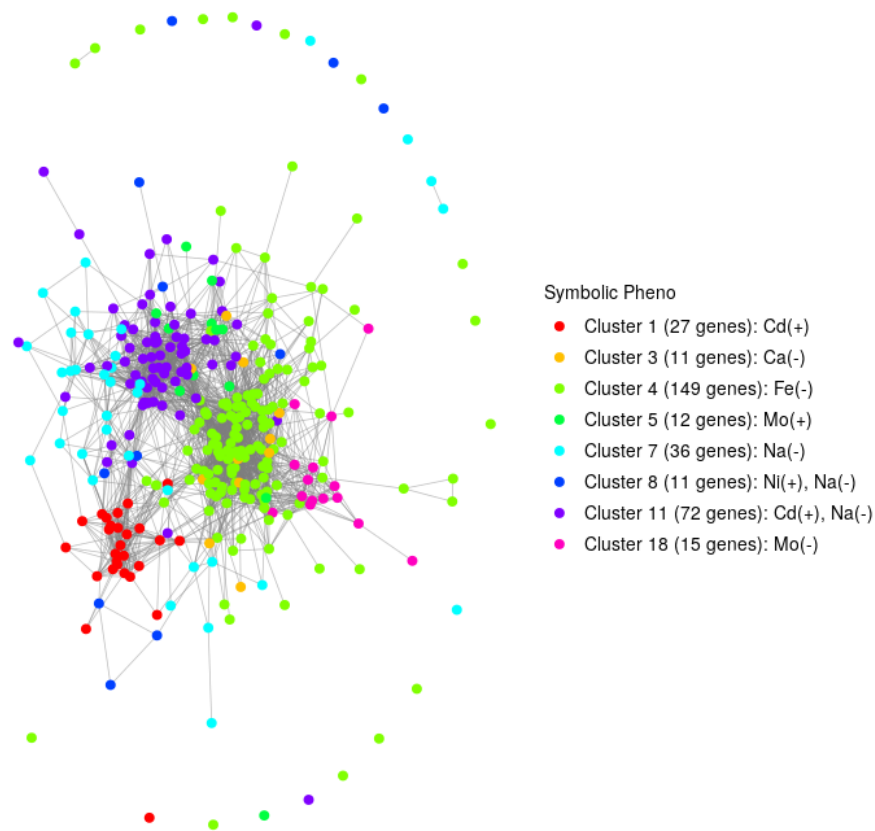


Figure 4: Network with Pearson correlation: impact and betweenness

## Ionflow: network and enrichment analysis for ionomics data

For comparison purposes, we use different similarity methods. Here we choose *Mahalanobis Cosine*:

```
net_2 <- GeneNetwork(data = dat,  
                      data_symb = dat_symb,  
                      min_clust_size = 10,  
                      thres_corr = 0.75,  
                      method_corr = "mahal_cosine")  
  
net_2$plot.pnet1
```



**Figure 5: Network with Mahalanobis Cosine**

```
net_2$plot.pnet2
```

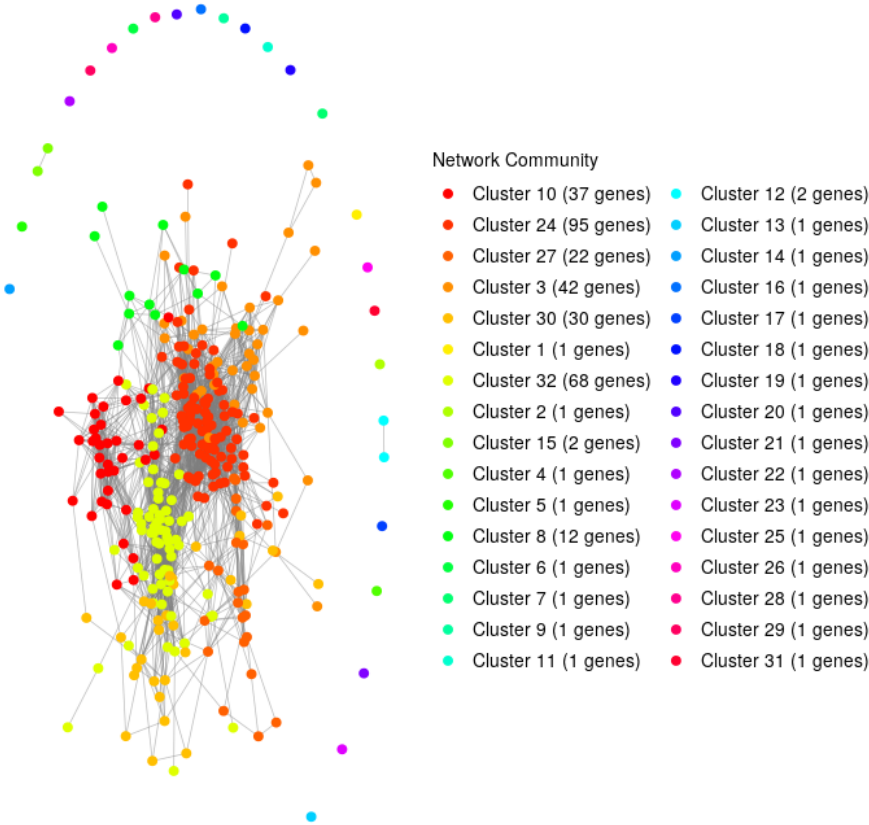


Figure 6: Network with Mahalanobis Cosine

## Enrichment analysis

The enrichment analysis is used for group data. The genes in groups are considered target gene sets while genes in the whole data set is the universal gene set. The group data can be results of the symbolic clustering or network community centres.

The Bioconductor R package [GOstats](#) is used for the enrichment analysis.

The network analysis returns a vertex attributes matrix:

```
head(net$net_node)
#>      Line          symb_pheno      comm_centre
#> 1 YAL009W Cluster 1 (27 genes): Cd(+) Cluster 8 (50 genes)
#> 2 YAL013W Cluster 3 (11 genes): Ca(-) Cluster 5 (162 genes)
#> 3 YAL014C Cluster 1 (27 genes): Cd(+) Cluster 8 (50 genes)
#> 4 YAL020C Cluster 4 (149 genes): Fe(-) Cluster 5 (162 genes)
#> 5 YAL021C Cluster 4 (149 genes): Fe(-) Cluster 5 (162 genes)
#> 6 YAL022C Cluster 1 (27 genes): Cd(+) Cluster 8 (50 genes)
```

The second and third columns are symbolic clustering and network community cluster, respectively.

If we perform enrichment analysis on the network community centre, the matrix should include the first column (gene IDs) and the third column.

The KEGG enrichment analysis, using p-values of 0.05 and genome wide annotation for Yeast, [org.Sc.sgd.db](#):

```
mat <- net$net_node[, c(1,3)]
kegg <- kegg_enrich(mat = mat, pval = 0.05, annot_pkg = "org.Sc.sgd.db")

#' kegg
kegg %>%
  kable(caption = 'KEGG enrichment analysis on network community centre',
        digits = 3, booktabs = T) %>%
  kable_styling(full_width = F, font_size = 10,
                latex_options = c("striped", "scale_down"))
```

**Table 6: KEGG enrichment analysis on network community centre**

comm_centre	KEGGID	Pvalue	Count	Size	Term
Cluster 2 (74 genes)	00400	0.025	2	2	Phenylalanine, tyrosine and tryptophan biosynthesis
Cluster 3 (39 genes)	00260	0.010	3	4	Glycine, serine and threonine metabolism
Cluster 3 (39 genes)	00290	0.021	2	2	Valine, leucine and isoleucine biosynthesis
Cluster 3 (39 genes)	00520	0.021	2	2	Amino sugar and nucleotide sugar metabolism
Cluster 6 (5 genes)	04011	0.006	2	4	MAPK signaling pathway - yeast
Cluster 8 (50 genes)	04111	0.044	2	5	Cell cycle - yeast

Note that there could be no results returned for KEGG enrichment analysis.

## Ionflow: network and enrichment analysis for ionomics data

The GO Terms enrichment analysis with ontology of *BP* (other two are *MF* and *CC*):

```
go <- go_enrich(mat = mat, pval = 0.05, ont = "BP", annot_pkg = "org.Sc.sgd.db")
#' go
dim(go)
#> [1] 45 7
go %>% head() %>%
  kable(caption = 'GO Terms enrichment analysis on network community centre',
        digits = 3, booktabs = T) %>%
  kable_styling(full_width = F, font_size = 10,
                latex_options = c("striped", "scale_down"))
```

**Table 7: GO Terms enrichment analysis on network community centre**

comm_centre	ID	Description	Pvalue	Count	CountUniverse	Ontology
Cluster 2 (74 genes)	GO:0007033	vacuole organization	0.0107	3	3	BP
Cluster 2 (74 genes)	GO:0005975	carbohydrate metabolic process	0.0372	8	19	BP
Cluster 2 (74 genes)	GO:0000291	nuclear-transcribed mRNA catabolic process, exonucleolytic	0.049	2	2	BP
Cluster 2 (74 genes)	GO:0002376	immune system process	0.049	2	2	BP
Cluster 2 (74 genes)	GO:0006952	defense response	0.049	2	2	BP
Cluster 2 (74 genes)	GO:0009073	aromatic amino acid family biosynthetic process	0.049	2	2	BP

We can also perform enrichment analysis on the symbolic clustering. To do so, use the first and second columns. KEGG enrichment analysis:

```
mat <- net$net_node[, c(1,2)]
kegg <- kegg_enrich(mat = mat, pval = 0.05, annot_pkg = "org.Sc.sgd.db")
kegg %>%
  kable(caption = 'KEGG enrichment analysis on symbolic clustering',
        digits = 3, booktabs = T) %>%
  kable_styling(full_width = F, font_size = 10,
                latex_options = c("striped", "scale_down"))
```

**Table 8: KEGG enrichment analysis on symbolic clustering**

symp_pheno	KEGGID	Pvalue	Count	Size	Term
Cluster 11 (72 genes): Cd(+), Na(-)	00400	0.021	2	2	Phenylalanine, tyrosine and tryptophan biosynthesis
Cluster 18 (15 genes): Mo(-)	01100	0.008	6	37	Metabolic pathways
Cluster 18 (15 genes): Mo(-)	00564	0.014	2	3	Glycerophospholipid metabolism
Cluster 7 (36 genes): Na(-)	00260	0.010	3	4	Glycine, serine and threonine metabolism
Cluster 7 (36 genes): Na(-)	00290	0.021	2	2	Valine, leucine and isoleucine biosynthesis
Cluster 7 (36 genes): Na(-)	00520	0.021	2	2	Amino sugar and nucleotide sugar metabolism

GO Terms enrichment analysis:

```
go <- go_enrich(mat = mat, pval = 0.05, ont = "BP", annot_pkg = "org.Sc.sgd.db")
dim(go)
#> [1] 71 7
go %>% head() %>%
  kable(caption = 'GO Terms enrichment analysis on symbolic clustering',
        digits = 3, booktabs = T) %>%
```

```
kable_styling(full_width = F, font_size = 10,
              latex_options = c("striped", "scale_down"))
```

Table 9: GO Terms enrichment analysis on symbolic clustering

symb_pheno	ID	Description	Pvalue	Count	CountUniverse	Ontology
Cluster 1 (27 genes): Cd(+)	GO:0051336	regulation of hydrolase activity	0.0019	4	8	BP
Cluster 1 (27 genes): Cd(+)	GO:0043085	positive regulation of catalytic activity	0.0032	4	9	BP
Cluster 1 (27 genes): Cd(+)	GO:0035303	regulation of dephosphorylation	0.0063	2	2	BP
Cluster 1 (27 genes): Cd(+)	GO:0046889	positive regulation of lipid biosynthetic process	0.0063	2	2	BP
Cluster 1 (27 genes): Cd(+)	GO:1903727	positive regulation of phospholipid metabolic process	0.0063	2	2	BP
Cluster 1 (27 genes): Cd(+)	GO:0044764	multi-organism cellular process	0.0131	3	7	BP



### Exploratory analysis

The exploratory analysis performs PCA and correlation analysis for ions in terms of genes. Note that this analysis treats ions as samples/replicates while genes are treated as variables/features. The exploratory analysis is initially employed at an early stage of the analysis.

For example, we apply it to the pre-processed data `dat` before any other analysis:

```
expl <- ExploratoryAnalysis(data = dat)
names(expl)
#> [1] "plot.pca"      "data.pca.load" "plot.corr"      "plot.corr.heat"
#> [5] "plot.heat"     "plot.net"
```

The PCA plot is:

```
expl$plot.pca
```

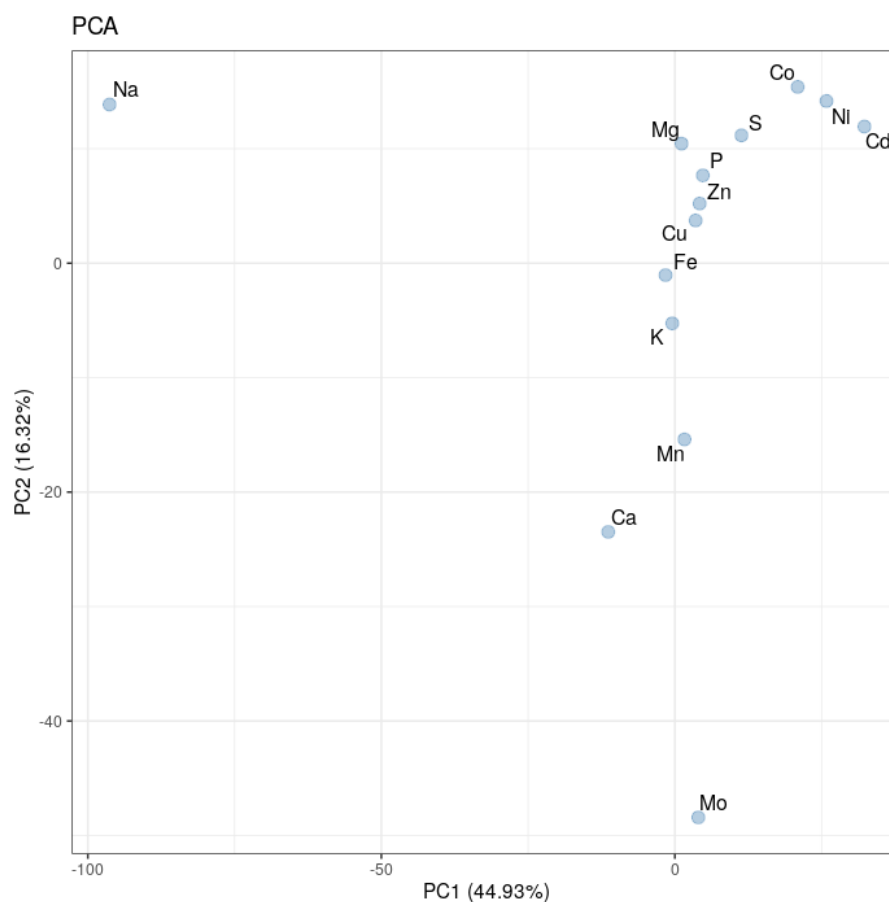


Figure 7: Ion PCA plot on pre-processed data

Ionflow: network and enrichment analysis for ionomics data

The Person correlation of ions are shown in correlation plot, heatmap and network plot:

```
expl$plot.corr
```

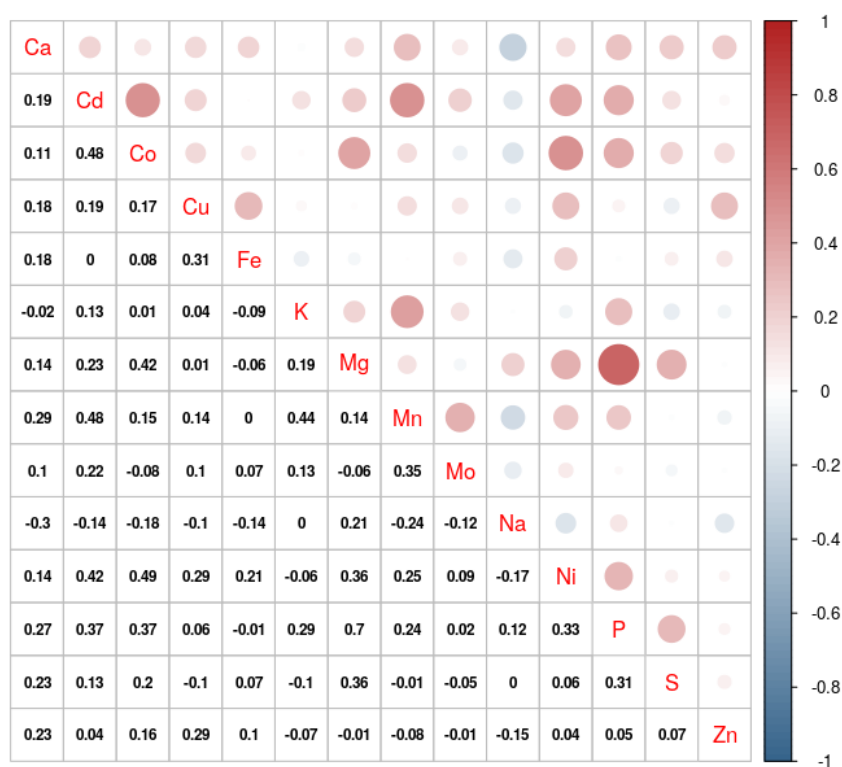


Figure 8: Ion correlation plots on pre-processed data

```
expl$plot.corr.heat
```

```
expl$plot.net
```

The correlation between ions and genes are shown in heatmap with dendrogram:

```
expl$plot.heat
```

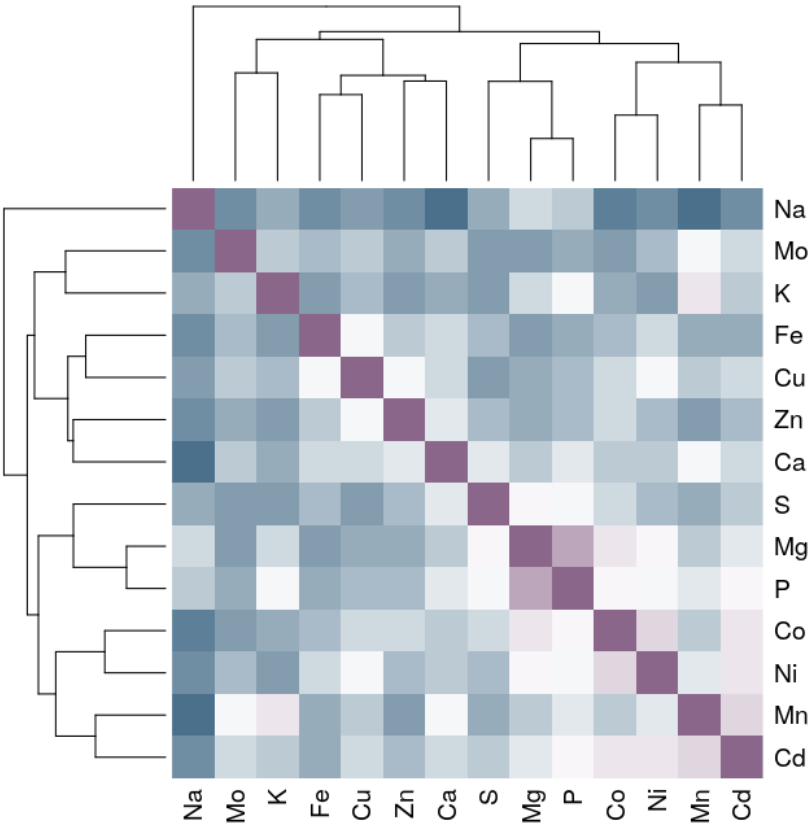


Figure 9: [ion correlation plots on pre-processed data](#)

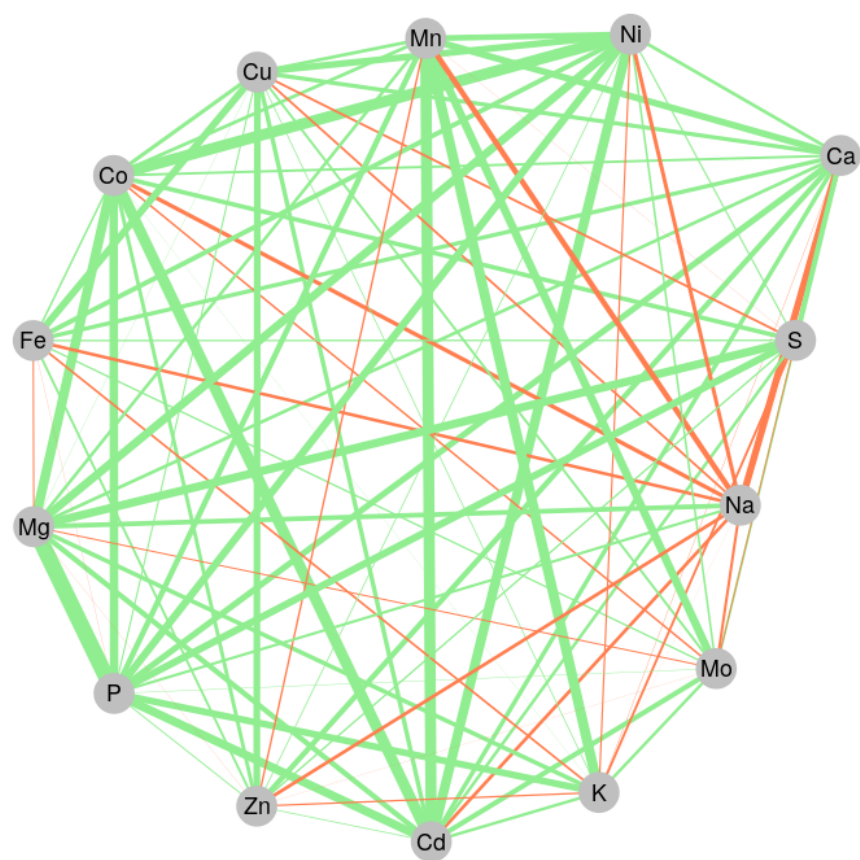


Figure 10: Ion correlation plots on pre-processed data

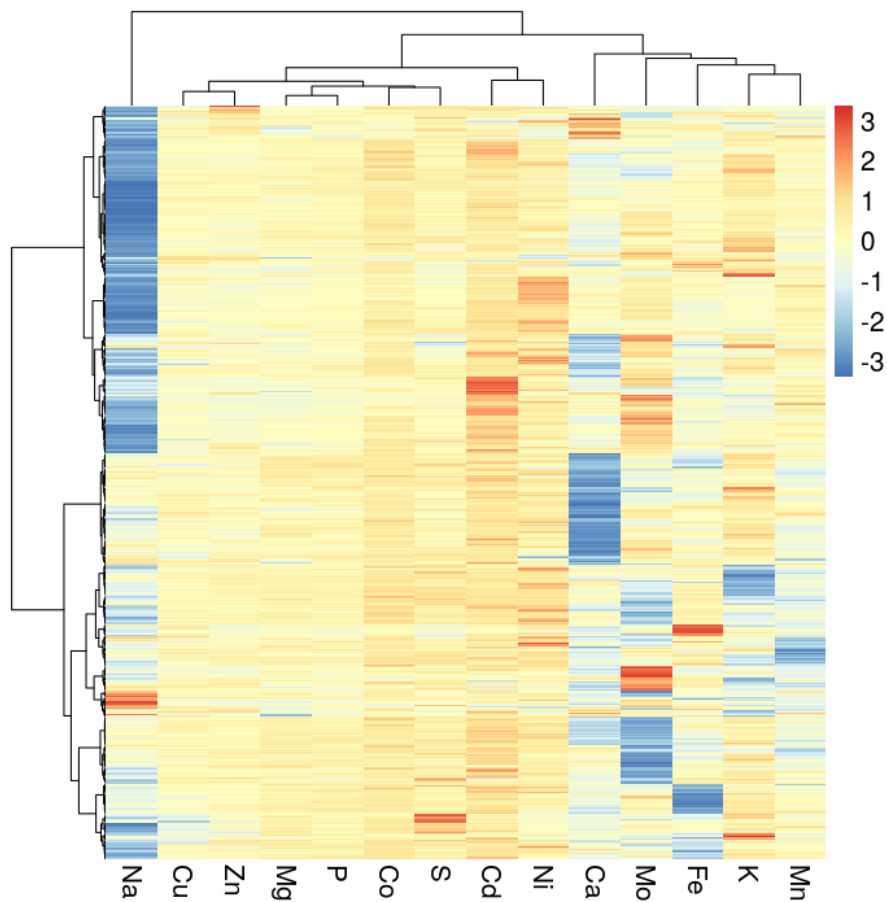


Figure 11: Correlation between ions and genes on pre-processed data

## Ionflow: network and enrichment analysis for ionomics data

The exploratory analysis can also be used at other stages of the analysis. Here for example after gene clustering analysis:

```
#' update data set with results of gene clustering
dat_clus <- dat[clust$idx, ]
dim(dat_clus)
#> [1] 333 15

expl.1 <- ExploratoryAnalysis(data = dat_clus)
```



Figure 12: Exploratory analysis after gene clustering

```
expl.1$plot.pca
```

```
expl.1$plot.net
```

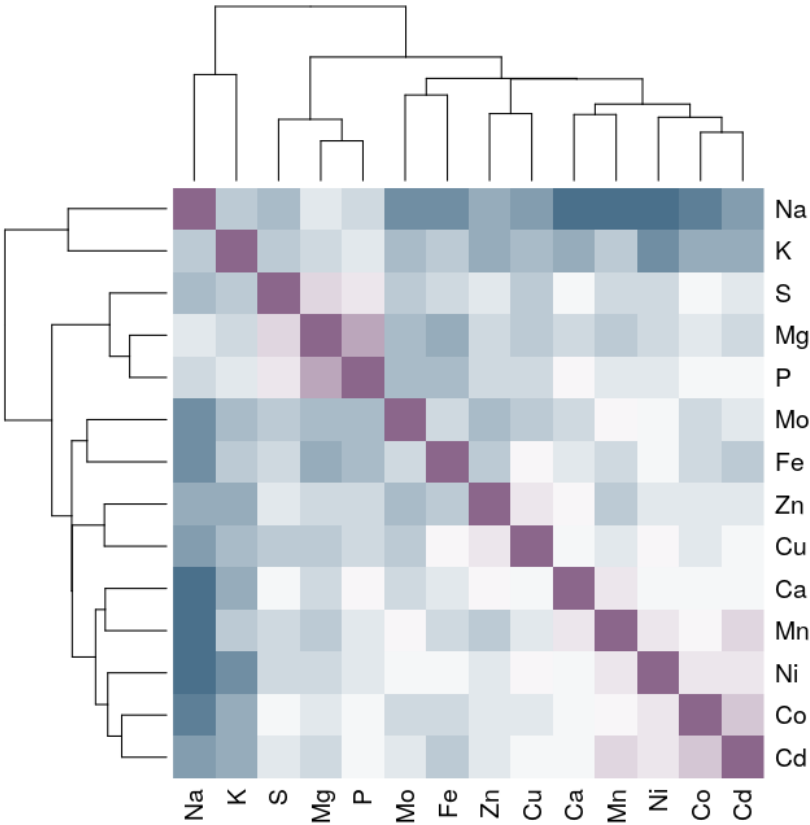


Figure 13: Exploratory analysis after gene clustering

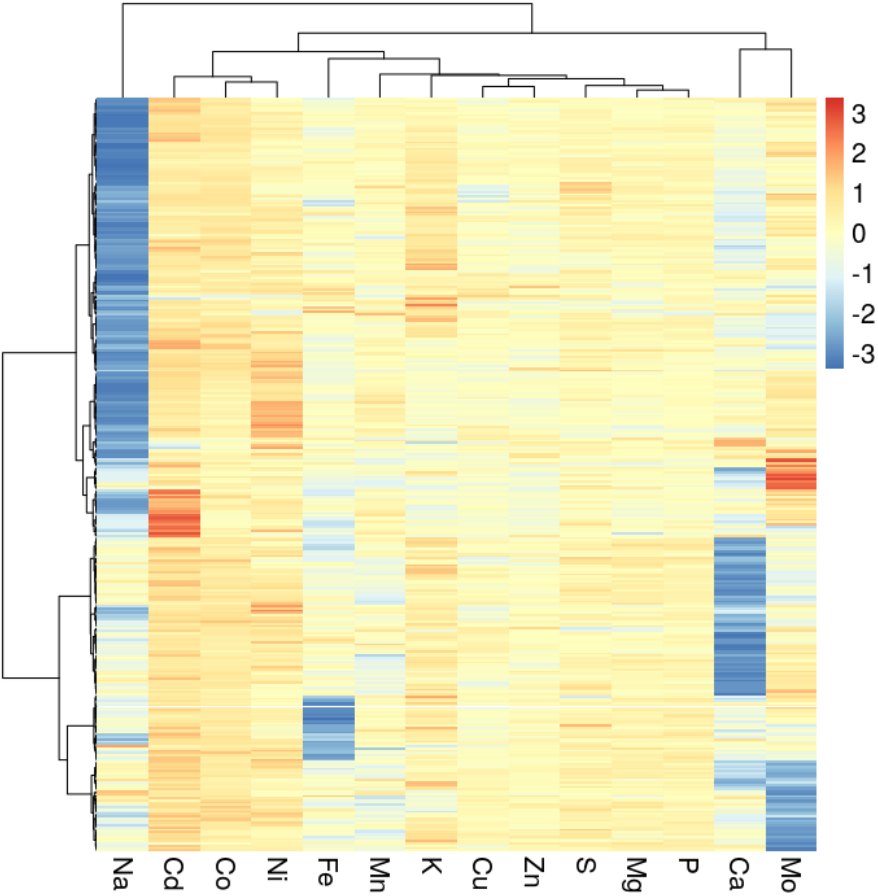


Figure 14: Exploratory analysis after gene clustering



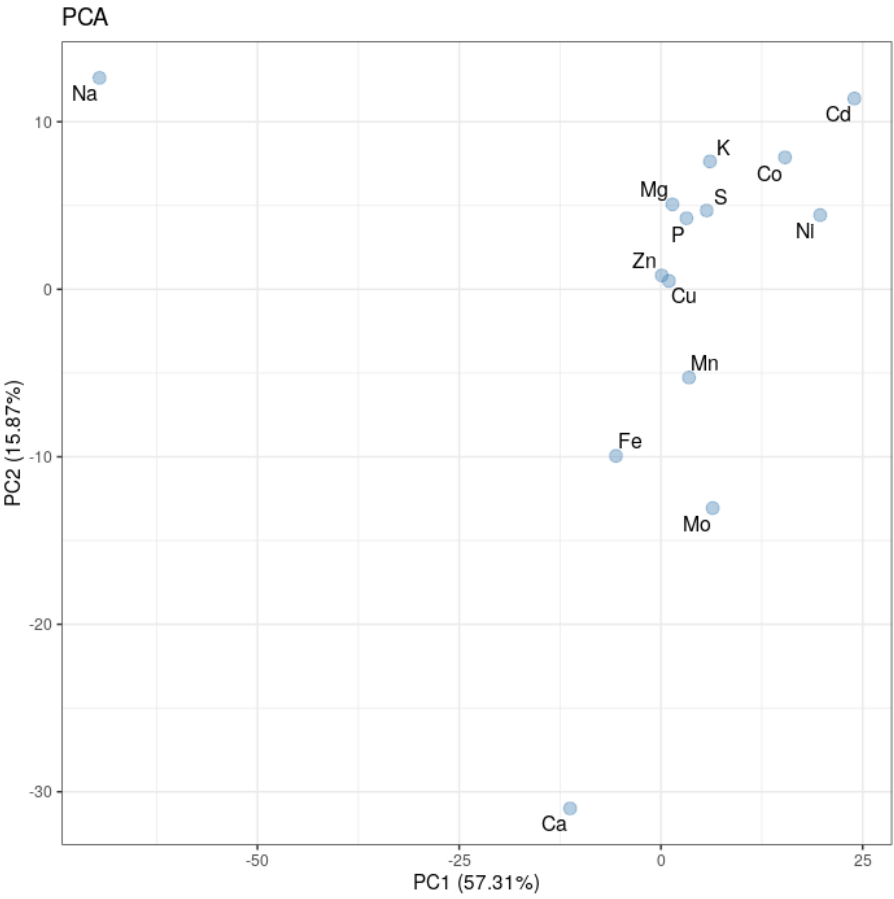


Figure 15: Exploratory analysis after gene clustering

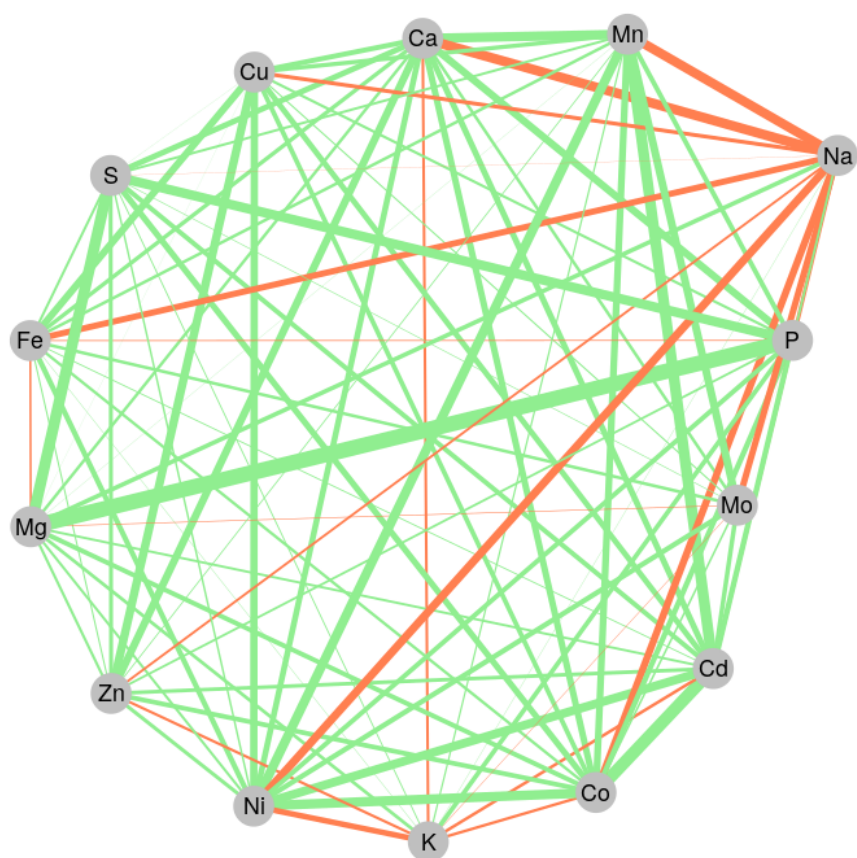


Figure 16: Exploratory analysis after gene clustering