

1. Open Terminal and navigate to the folder containing the MetICA Code

```
1 cd directory/with/code
```

2. Create a directory named 'Results' lying in the same parent directory as the code. This folder will contain files produced by the code

```
1 cd ..  
2 mkdir Results
```

3. Open R Console and execute the following steps to load all libraries and scripts and to parse and preprocess the data.

**Note:** I assume that the data are saved in a csv-file, where the first column contains the ID and the following columns represent samples while the rows represent peaks. Each entry corresponds to the intensity of a peak in the given sample. If your data is not structured in this form yet, please re-structure it before executing these steps

```
1 source('MetICA_load_all.R')  
2  
3 example_data = read.table('example.csv', sep=',', dec='.',  
4 header=TRUE, check.names=FALSE)  
5 peaks = data.matrix(example_data[, 2:ncol(example_data)])  
6  
7 #Log-scaling of peaks  
8 peaks = log(peaks + min(peaks) + 1)  
9  
10 #Taking only peaks appearing in more than one sample...  
11 peaks = peaks[apply(peaks>0, MARGIN=1, sum) > 1, ]  
12  
13 #and having different values for at least 2 samples  
14 ranges = apply(peaks, MARGIN=1, range)  
15 peaks = peaks[ranges[, 1] != ranges[, 2], ]
```

4. Start generating independent components. (The terms 'independent components' and 'sources' are used interchangeably here.) This will output a file named `source_list_<nbcomp>.rda`, where `<nbcomp>` is replaced by the number of independent components.

**Note:** The second parameter of the function specifies the total variance kept while the fourth parameter specifies the number of iterations. Both are important for the expected time of computation. In this example, only 30% of variance is kept to make things run faster. The total variance will also directly determine the number of principal components which is equal to the number of sources.

```
1 M1 = MetICA_source_generator(peaks, 0.3, 'gaussian', 800)
```

5. Cluster the generated sources: The cluster generator evaluates the optimal number of clusters and computes centrotypes for this number. It produces the following output files in the 'Results' folder:

- `nbclust_eval_<nbcomp>_<nbclust>.csv`, where `<nbcomp>` is the number of components and `<nbclust>` is the determined optimal number of clusters. The file contains average silhouette scores for different numbers of clusters

- `center_<nbclust>.rda`, which contains the resulting centrotypes

This will also create several auxiliary files in the 'Dist' folder, storing dissimilarity matrices. Before running the code again, all files from this folder should be deleted.

```
1 M2 = MetICA_cluster_generator(M1$S, 'spearman', M1$nbcomp)
```

6. (Optional) An evaluation of statistical reliability of the resulting centrotypes can be done on bootstrapped datasets. This will produce the following outputs in 'Results'

- `IC_notes_summary_<nbclust>.csv`, storing scores for each bootstrapped dataset (rows) and centrotype (columns)
- `IC_ranking_<nbclust>.png`, showing the corresponding ranking of centrotypes in a boxplot

**Warning:** This might take a long time even for small numbers of components.

```
1 M3 = MetICA_bootstrap(peaks, 5, M2$centers, M1$nbcomp,
2                       'gaussian', 100, 50)
```