

# Package ‘MIA’

**Version** 1.0

**Date** 2016-01-14

**Title** Matrix Integration Analysis

**Author** Jinyu Chen <chjy@amss.ac.cn>, Shihua Zhang <zsh@amss.ac.cn>

**Depends** MATLAB ( $\geq$  R2013a)

**Description** Detecting multi-dimensional modules (md-modules) in diverse genomics data as well as molecular network data using the methods in MIA package.

**URL** <http://page.amss.ac.cn/shihua.zhang/software.html>

Here, we provide a guide for MIA package. It describes all the MATLAB functions in MIA in detail. For each method, these MATLAB functions mainly perform the tasks including realizing a specific algorithm, drawing figures and outputting text files about the identified md-modules.

# 1 jNMF

**jNMF** (joint Non-negative Matrix Factorization) enables users to simultaneously factor two or more types of genomic data sharing the same set of samples. Here, we adopt multiplicative update algorithm to solve the following problem:

$$\min_{W, H_i} \sum_{i=1}^N \|X_i - WH_i\|_F^2, \quad \text{s.t. } W \geq 0, H_i \geq 0, i = 1, \dots, N.$$

## 1.1 Algorithm

---

Run_jNMF	<i>The main function for jNMF.</i>
----------	------------------------------------

---

**Description**

This is the main function for jNMF, which integrates all the related functions to achieve it.

**Usage**

```
Run_jNMF(Input);
```

**Arguments**

Input	A structure variable. The details about its construction can be found in Section 5.
-------	---

**Output**

It saves all the results in the directory `./MIA/jNMF/jNMF_Results/`.

<i>jNMF_Results.mat</i>	The essential computed variables.
<i>jNMF_RunRecords.txt</i>	The updated values of the objective function in each round of running.
<i>jNMF_Results.txt</i>	The numbers and indexes of features separated in comma. Each row records the information for each identified md-module.
Several folders	Each folder contains the lists of all the identified md-modules for one type of features as shown in Figure 2a and Figure 3a.
Several figures	As shown in Figure 2b and Figure 4.

---

jNMF_PrepData	<i>Preprocess the input data.</i>
---------------	-----------------------------------

---

**Description**

This function is used to preprocess the input data to ensure its non-negativity.

**Usage**

```
[X, isdouble] = jNMF_PrepData(OX);
```

**Arguments**

OX	A matrix.
----	-----------

**Output**

X	A non-negative matrix, the transformation of data OX.
isdouble	A binary variable (1 indicates changes have been made; 0 is for no change, that is, $X = OX$ ).

---

jNMF_comodule	<i>Obtain the md-modules.</i>
---------------	-------------------------------

---

## Description

This function outputs the optimal factorization results through running jNMF for multiple times, then identify md-modules based on the factorized matrices  $W$ ,  $H_i$  ( $i = 1, 2, \dots$ ).

## Usage

```
[W, H, Comodule, params] = jNMF_comodule(Input, params);
```

## Arguments

Input	A structure variable including two components:
Input.data	A non-negative matrix combing $N$ data blocks sequentially to be factorized, such as $Input.data = [X_1, X_2, \dots, X_N]$ .
Input.XBlockInd	A matrix of size $N \times 2$ . The two elements of the $i$ th row give the start and end column indexes in $Input.data$ for data matrix $X_i$ ( $i = 1, \dots, N$ ).
params	A structure variable including six components:
params.isdouble	A vector of size $N \times 1$ . The $i$ th element indicates whether the $i$ th data matrix is transformed to ensure its non-negativity (0 for no change, and 1 for change).
params.K	The number of md-modules users prefer to identify.
params.nloop	The repeating times of jNMF. To ensure the robust of this method, this function repeats the algorithm for 'params.nloop' times.
params.maxiter	The maximal number of iterations for this algorithm.
params.tol	The precision for convergence of algorithm.
params.thrd_module	A positive vector of size $1 \times (N + 1)$ . Thresholds for selecting features in $N$ data blocks. The first one is to select samples.

## Output

W, H	Factorization results such that $Input.data \approx WH$ . $W$ is the basis matrix of size $m \times K$ and $H$ is the weight matrix of size $K \times n$ , where $K = params.K$ .
Comodule	Identified md-modules recorded in a $(K \times (N + 1))$ cell array. $Comodule\{i,j\}$ records selected feature indexes of the $j$ th type of variables in the $i$ th identified md-module. The first column is for selected samples.
params	Compared to 'params' as input, there are something new added in it, including
params.records	A $(nloop \times 1)$ cell array. $params.records\{i\}$ is a $(iter \times (N + 1))$ vector, where each row records the values of all the terms in the objective function and the sum of them in each iteration.
params.iterNumList	A $(nloop \times 1)$ vector, where each element is the number of iterations for each round of running.

---

jNMF_algorithm	<i>jNMF algorithm.</i>
----------------	------------------------

---

## Description

This is jNMF algorithm.

## Usage

```
[W, H, TerminalObj, iter] = jNMF_algorithm(X, XInd, params);
```

## Arguments

X	A non-negative input matrix of size $m \times n$ combing $N$ data matrices sequentially to be factorized (e.g., $X = [X_1, X_2, \dots, X_N]$ ).
XInd	A matrix of size $N \times 2$ . The two elements of the $i$ th row give the start

params and end column indexes in  $X$  for data matrix  $X_i$  ( $i = 1, \dots, N$ ).  
 A structure variable including *params.K*, *params.maxiter*, *params.tol* defined the same as that in the function ‘*jNMF\_comodule*’.

### Output

W, H Factorization results like those in function ‘*jNMF\_comodule*’.  
 iter The number of iterations when the algorithm stops.  
 TerminalObj A ( $iter \times (N + 1)$ ) matrix in which each row records the values of all the terms in the objective function and the sum of them in each iteration.

*jNMF\_module*                      *Identify md-modules from factorized matrices.*

### Description

Based on the factorized matrix  $W$  or  $H_i$ , identify module members for each type of features.

### Usage

module = *jNMF\_module*(H, t, isdouble);

### Arguments

H A ( $K \times m$ ) non-negative matrix used for module identification.  
 t A threshold value for selecting features.  
 isdouble A binary variable (1 is for the number of features in matrix  $H$  is double than that of the original ones, and 0 is for no change).

### Output

module A ( $K \times 1$ ) cell array. module*i*,1 contains the feature indexes of the *i*th module.

## 1.2 Output figures

*jNMF\_plot\_X*                      *Provide the heatmaps of the original input matrices.*

### Description

Draw the heatmaps of the original input matrices ( $X_1, X_2, \dots, X_N$ ).

### Usage

*jNMF\_plot\_X*(X, XInd, fig, figure\_title, colormap\_type);

### Arguments

X The input matrix combing  $N$  data matrices to be factorized (e.g.,  $X = [X_1, X_2, \dots, X_N]$ ).  
 XInd A matrix of size  $N \times 2$ . The two elements of the *i*th row give the start and end column indexes in  $X$  for data matrix  $X_i$  ( $i = 1, \dots, N$ ).  
 fig A positive integer for figure index.  
 figure\_title A string for the title of figure.  
 colormap\_type A string for the colormap of heatmaps. Options includes ‘blue-yellow’, ‘green-red’, ‘yellow’, ‘blue-white-red’, ‘default’.

### Output

The heatmaps for all the input data matrices as shown in Figure 4.

*jNMF\_plot\_results*                      *Show the heatmaps of certain identified md-modules.*

---

**Description**

Show the heatmaps of a selected identified md-module (circled in yellow lines).

**Usage**

```
jNMF_plot_results(X, XInd, FeatureType, fig, figure_title, colormap_type, vectorForRank);
```

**Arguments**

*X, XInd, fig, figure\_title, colormap\_type*

They have the same definitions with those in function *jNMF\_plot\_X*.

FeatureType

A  $(1 \times N)$  cell array. *FeatureType* $\{i\}$  records the name of *i*th type of features (e.g., *FeatureType* = {'Gene expression', 'microRNA expression', 'CNV'}).

vectorForRank

A structure variable containing four components:

vectorForRank.w, vectorForRank.h

Two vectors for the selected md-module. For example, if one wants to demonstrate the *i*th md-module, *vectorForRank.w* is the *i*th column of basis matrix *W* and *vectorForRank.h* is the *i*th row of weight matrix *H*.

vectorForRank.comodule

A  $(1 \times (N + 1))$  cell array. *vectorForRank.comodule* $\{i\}$  records the *i*th feature indexes of the selected md-module.

vectorForRank.hInd

Similar with the input variable 'XInd' to 'X', it records the indexes for '*vectorForRank.h*'.

**Output**

The heatmaps for a selected md-module to demonstrate its patterns as shown in Figure 4.

---

jNMF\_plot\_correlation

*Demonstrate the correlations between the original data and reconstructed data.*

---

**Description**

Demonstrate the correlations between the original data  $X_i$  and reconstructed data  $newX_i = WH_i$  using boxplots.

**Usage**

```
corrMat = jNMF_plot_correlation(X, newX, XInd, newXInd, fig, figure_title);
```

**Arguments**

*X, XInd, fig, figure\_title*

They have the same definitions with those in function '*jNMF\_plot\_X*'.

newX

The reconstructed matrix combing  $N$  data matrices sequentially, that is,  $newX = [WH_1, WH_2, \dots, WH_N]$ .

newXInd

Similar with the input variable 'XInd' to 'X', it records the indexes for matrix *newX*.

**Output**

corrMat

A matrix of size  $m \times N$ , where  $m$  is the number of samples and  $N$  is the input data matrices, respectively. *corrMat* $(i,j)$  records the correlation between the *i*th rows of the original data ( $X_j$ ) and reconstructed data ( $WH_j$ ).

Boxplots for input data matrices  $X_1, X_2, \dots, X_N$  as shown in Figure 4.

---

jNMF\_plot\_distribution

*Demonstrate the module size distributions.*

---

## Description

This function provides the histograms for the size distributions of  $(N + 1)$  types of features in the identified md-modules.

## Usage

```
jNMF_plot_distribution(nSample, XInd, Comodule, FeatureType, fig, figure_title);
```

## Arguments

*XInd, FeatureType, fig, figure\_title*

They are the same as those in function ‘*jNMF\_plot\_results*’.

nSample                   The number of samples.

Comodule                  It is the same as that in function ‘*jNMF\_comodule*’.

## Output

Histograms for the size distributions of  $(N + 1)$  types of features in the identified md-modules as shown in Figure 4.

## 1.3 Output into text files

---

Index2LabelForModuleContent	<i>Output the identified md-modules into text files.</i>
-----------------------------	--

---

## Description

Output a number of text files, each of which records the selected feature names.

## Usage

```
Index2LabelForModuleContent(ModuleIndex, FeatureLabel, TypeName, ResultsFile);
```

## Arguments

ModuleIndex              A  $(K \times 1)$  cell array. *ModuleIndex* $\{i\}$  records the indexes of the *i*th identified module for one type of features.

FeatureLabel             A  $(n \times 1)$  cell array recording all the names of one type of features.

TypeName                 A string for the feature type (e.g., *TypeName* = ‘*Gene expression*’).

ResultsFile              A string for the folder name to save these lists (e.g., *TypeName* = ‘*Gene expression*’, *ResultsFile* = ‘*jNMF\_results*’, all the identified gene lists are saved in the directory: ‘*./jNMF\_results/GeneLists/*’).

## Output

A number of text files (e.g., *GeneList\_1.txt*, *GeneList\_2.txt*, ...)

---

OutputModule2TXT	<i>Output the feature indexes of the identified md-modules.</i>
------------------	---

---

## Description

Output the feature indexes of the identified md-modules into a text file.

## Usage

```
OutputModule2TXT(Comodule, FeatureType, ResultsFile);
```

## Arguments

Comodule                 It is the same as that in function ‘*jNMF\_comodule*’.

ResultsFile              A string for the name of this text file.





## Arguments

Input	A structure variable including three components:
Input.data	A matrix combing two data blocks sequentially to be factorized (e.g., $Input.data = [X_1, X_2]$ ).
Input.XBlockInd	A matrix of size $2 \times 2$ . The two elements of the $i$ th row give the start and end column indexes in $Input.data$ for data matrix $X_i$ ( $i = 1, 2$ ).
Input.netAdj	An adjacency matrix for the relationships between the features in $Input.data$ , that is, $Input.netAdj = [A_{11}, A_{12}; A_{12}^T, A_{22}]$ .

## Output

newInput	A structure variable including
newInput.data	A non-negative matrix which is transformed from $Input.data$ .
newInput.XBlockInd	A matrix of size $2 \times 2$ which is similar with $Input.XBlockInd$ .
newInput.netAdj	An adjacency matrix transformed from $Input.netAdj$ .
isdouble	A binary variable (1 indicates changes have been made, and 0 is for no change, that is, $Input = newInput$ ).

---

SNMNMF\_comodule

*Obtain the md-modules.*

---

## Description

This function computes the optimal factorization results through running SNMNMF for multiple times, then identify the md-modules based on the factorized matrices  $W, H_i$ .

## Usage

```
[W, H1, H2, Comodule, params] = SNMNMF_comodule(Input, params);
```

## Arguments

Input	A structure variable defined as that in function 'SNMNMF_PrepData', but the $Input.data$ is a non-negative matrix.
params	A structure variable, except for $params.isdouble$ , $params.K$ , $params.nloop$ , $params.maxiter$ , $params.thrd\_module$ , defined similar as those in the function 'jNMF_comodule' described above. There are five specific components, including
params.thrXr, params.thrXc	parameters referring to the basis matrix $W$ -related term, the weight matrices $H_i$ -related terms in the objective function to limit the growth of $W$ ,
params.thrNet11, params.thrNet12, params.thrNet22	make $H_i$ sparse, respectively. parameters referring to the network $A_{11}, A_{12}, A_{22}$ related constraints in the objective function.

## Output

W, H1, H2	Factorization results such that $Input.data \approx [WH_1, WH_2]$ . $W$ is the basis matrix of size $m \times K$ , $H_1$ is the weight matrix of size $K \times n_1$ and $H_2$ is the weight matrix of size $K \times n_2$ , where $K = params.K$ .
Comodule	Identified md-modules recorded in a $(K \times 3)$ cell array, which has the same definition as that in function 'jNMF_comodule'.
params	Compared to 'params' as input, there are something new added in it, including $params.records$ and $params.iterNumList$ which are similar with those in the function 'jNMF_comodule', where $N = 2$ .

---

SNMNMF_algorithm	<i>SNMNMF algorithm.</i>
------------------	--------------------------

---

**Description**

This implements the SNMNMF algorithm.

**Usage**

[W, H1, H2, TerminalObj, iter] = SNMNMF\_algorithm(X1, X2, A11, A12, A22, params);

**Arguments**

X1, X2	The non-negative input matrices.
A11, A12, A22	Adjacency matrices for the relationships within and between the features in $X_1$ and $X_2$ .
params	A structure variable including <i>params.isdouble</i> , <i>params.K</i> , <i>params.thrd_module</i> , <i>params.thrXr</i> , <i>params.thrXc</i> , <i>params.thrNet11</i> , <i>params.thrNet12</i> , <i>params.thrNet22</i> . They are the same as those in function ‘SNMNMF_comodule’.

**Output**

W, H1, H2	Factorization results like those in function ‘SNMNMF_comodule’.
TerminalObj	A ( $iter \times 3$ ) matrix each row of which records the values of all the terms in the objective function and the sum of them in each iteration.

---

SNMNMF_module	<i>Identify the md-modules from factorized matrix.</i>
---------------	--

---

**Description**

Based on the factorized matrix  $W$ ,  $H_1$  or  $H_2$ , identify the md-module members for each type of features.

**Usage**

module = SNMNMF\_module(H, t, isdouble);

**Arguments**

H, t, isdouble	They are defined as those in function ‘jNMF_module’.
----------------	--

**Output**

module	Defined as that in function ‘jNMF_module’.
--------	--

**2.2 Output figures**

---

SNMNMF_plot_X	<i>Provide the heatmaps for the original input matrices.</i>
---------------	--

---

**Description**

Draw the heatmaps for the original input matrices ( $X_1$  and  $X_2$ ). The arguments and outputs of this function are similar with those in the function ‘jNMF\_plot\_X’ of jNMF with  $N = 2$ .

**Usage**

SNMNMF\_plot\_X(X, XInd, fig, figure\_title, colormap\_type);

---

SNMNMF_plot_results	<i>Show the heatmaps of a identified md-module.</i>
---------------------	---

---

**Description**

Show the heatmaps of a selected identified md-module (circled in yellow lines). The arguments and outputs of this function are similar with those in the function 'jNMF\_plot\_result' with  $N = 2$ .

#### Usage

SNMNMF\_plot\_results(X, XInd, FeatureType, fig, figure\_title, colormap\_type, vectorForRank);

---

SNMNMF_plot_correlation	<i>Demonstrate the correlations between the original data and the reconstructed one.</i>
-------------------------	--

---

#### Description

Demonstrate the correlations between the original data  $X_i$  and the reconstructed one  $newX_i = WH_i$  using boxplot ( $i = 1, 2$ ). The arguments and outputs of this function are similar with those in the function 'jNMF\_plot\_correlation' with  $N = 2$ .

#### Usage

corrMat = SNMNMF\_plot\_correlation(X, newX, XInd, newXInd, fig, figure\_title);

---

SNMNMF_plot_distribution	<i>Demonstrate the module size distributions.</i>
--------------------------	---

---

#### Description

This function provides the histograms for the size distributions of two types of features in the identified md-modules. The arguments and outputs of this function are similar with those in the function 'jNMF\_plot\_distribution' with  $N = 2$ .

#### Usage

SNMNMF\_plot\_distribution(nSample, XInd, Comodule, FeatureType, fig, figure\_title);

### 2.3 Output text files

---

Index2LabelForModuleContent	<i>Output the identified md-modules into text files.</i>
-----------------------------	--

---

#### Description

This function is the same as that described in jNMF.

---

OutputModule2TXT	<i>Output the feature indexes of the identified md-modules.</i>
------------------	---

---

#### Description

This function is the same as that in jNMF.

## 3 sMBPLS

sMBPLS (sparse Multi-Block Partial Least Square) extends the standard PLS method to discover associations between multiple input matrices ( $X_1, X_2, \dots, X_N; N \geq 1, X_i \in \mathbb{R}^{p \times n_i}$  and a response matrix ( $Y \in \mathbb{R}^{p \times m}$ ) in a sparse manner. It identifies md-modules in which a subset of heterogeneous input features jointly explain a subset of the response variables. This problem is defined as,

$$\begin{aligned} & \max_{w_i, q, b_i} \text{cov}(t, u) - \sum_{i=1}^N \lambda_i \|w_i\|_1 - \tau \|q\|_1 \\ & \text{with } t_i = X_i w_i, u = Yq, t = \sum_{i=1}^N b_i t_i \\ & \text{s.t. } \|w_i\|_2^2 = 1, \|q\|_2^2 = 1, \sum_{i=1}^N b_i^2 = 1. \end{aligned}$$

### 3.1 Algorithm

---

<code>Run_sMBPLS</code>	<i>The main function for sMBPLS.</i>
-------------------------	--------------------------------------

---

**Description**

This is the main function for sMBPLS which integrates all the related functions to achieve it.

**Usage**

```
Run_sMBPLS(Input);
```

**Arguments**

Input	A structure variable (Section 5).
-------	-----------------------------------

**Output**

It saves all the results in the directory `./MIA/sMBPLS/sMBPLS_Results/`, including `sMBPLS_Results.mat`, `sMBPLS_RunRecords.txt`, `sMBPLS_Results.txt`, several folders and figures, which record the similar contents with those in function `'Run_jNMF'` of jNMF as shown in Figure 2, Figure 3 and Figure 4.

---

<code>meanc</code>	<i>Preprocess the input matrices.</i>
--------------------	---------------------------------------

---

**Description**

This function centers the input data across the samples to enable the mean of each column of the input matrix to be zero.

**Usage**

```
[Zm,mz] = meanc(Z);
```

**Arguments**

Z	A matrix.
---	-----------

**Output**

Zm	The centered matrix.
mz	A row vector. $mz(i)$ is the mean of $Z(:,i)$ .

---

<code>sMBPLS_comodule</code>	<i>Obtain the md-modules.</i>
------------------------------	-------------------------------

---

**Description**

This function returns all the md-modules by running SMBPLS for multiple times.

**Usage**

```
[nfactor, W, Q, WT, WU, TT, UU, sT, sU, XX, YY, Comodule, params] =  
sMBPLS_comodule(X, Y, XInd, YInd, FeatureType, params);
```

## Arguments

X	The input matrix of size $(p \times n)$ combining $N$ data blocks (e.g., $X = [X_1, \dots, X_N]$ ).
Y	The response matrix of size $(p \times m)$ combining $M$ data blocks (e.g., $Y = [Y_1, \dots, Y_M]$ ). Generally, $M = 1$ .
XInd	A matrix of size $N \times 2$ . The two elements of the $i$ th row give the start and end column indexes in $X$ for data matrix $X_i$ ( $i = 1, \dots, N$ ).
YInd	A matrix of size $M \times 2$ . It is for $Y$ defined similarly with $XInd$ .
FeatureType	A $(1 \times (N + M))$ cell array. $FeatureType\{i\}$ records the name of $i$ th type of features (e.g., $FeatureType = \{ 'microRNA\ expression', 'CNV', 'DNA\ Methylation', 'Gene\ expression' \}$ ).
params	A structure variable including
params.nfactor	A pre-defined number of identified md-modules.
params.nfold	A positive integer which is the number of folds used for cross-validation (CV) procedures. Generally, we set $params.nfold = 5$ or $10$ .
params.maxiter	The maximal number of iterations for sMBPLS algorithm.
params.tol	The precision for the convergence of sMBPLS algorithm.
params.param	A $(l \times 1)$ cell array restoring all the combinations of parameters to be selected. $params.param\{i\}$ contains one group of parameters used in the algorithm, including $thrXc$ and $thrYc$ (controlling the number of selected features of $X$ and $Y$ in the identified md-modules), $thrXr$ and $thrYr$ ( $thrXr = thrYr$ , controlling the number of selected samples in the identified md-modules), $maxiter$ ( $maxiter = params.maxiter$ ), $tol$ ( $tol = params.tol$ ).

## Output

nfactor	The number of identified md-modules ( $nfactor \leq params.nfactor$ ).
W	A $(n \times nfactor)$ matrix, $W(:,i)$ is the weight vector for $X$ in $i$ th md-module.
Q	A $(m \times nfactor)$ matrix, $Q(:,i)$ is the weight vector for $Y$ in $i$ th md-module.
WT	A $(N \times nfactor)$ matrix, $WT(:,i)$ is the super weight vector for $X$ in $i$ th md-module.
WU	A $(M \times nfactor)$ matrix, $WU(:,i)$ is the super weight vector for $Y$ in $i$ th md-module.
TT	A $(p \times (N \times nfactor))$ matrix, $TT(:, N \times (i - 1) + j)$ is the score vector for $X_j$ in $i$ th md-module ( $i = 1, \dots, nfactor; j = 1, \dots, N$ ).
UU	A $(p \times (M \times nfactor))$ matrix, $UU(:, M \times (i - 1) + j)$ is the score vector for $Y_j$ in $i$ th md-module ( $i = 1, \dots, nfactor; j = 1, \dots, M$ ). Generally, $M = 1$ .
sT	A $(p \times nfactor)$ matrix, $sT(:, i)$ is the super score for $X$ in $i$ th md-module ( $i = 1, \dots, nfactor$ ).
sU	A $(p \times nfactor)$ matrix, $sU(:, i)$ is the super score for $Y$ in $i$ th md-module ( $i = 1, \dots, nfactor$ ).
XX, YY	The new matrix after removing the signals of identified md-modules from $X$ and $Y$ .
Comodule	The md-modules are recorded in a $(nfactor \times (M + N + 1))$ cell array which is similar to that in the function 'jNMF_comodule'.
params	Compared to 'params' as input, there are something new added in it, including
params.iterNumList	A $(nfactor \times 1)$ vector, where each element is the number of iterations for each round of running.

params.records	A ( $nfactor \times 1$ ) cell array. $params.records\{i\}$ is a ( $iter \times (M + N + 2)$ ) matrix where each row records the values of all the terms in the objective function and the sum of them in each iteration when identify the $i$ th md-module, and $iter = params.iterNumList(i)$ , $i = 1, \dots, nfactor$ .
params.randRowPartitions	A ( $params.nfold \times 1$ ) cell array. For n-fold CV procedure, $p$ samples are partitioned into $params.nfold$ groups equally. $params.randRowPartitions\{i\}$ records the sample indexes in the $i$ th group ( $i = 1, \dots, params.nfold$ ).
params.cv_scores	A matrix where $params.cv\_scores(i, j)$ is the cv score for the $i$ th group of parameters when identify the $j$ th md-module.
params.paramsidx_used	A ( $nfactor \times 1$ ) vector. $params.paramsidx\_used(i)$ is the index of selected parameters for identifying the $i$ th md-module.

---

sMBPLS_algorithm	<i>sMBPLS algorithm.</i>
------------------	--------------------------

---

### Description

This implements the sMBPLS algorithm.

### Usage

[success, w, q, b, a, T, U, t, u, XX, YY, TerminalObj, iter]  
= sMBPLS\_algorithm(X, Y, XInd, YInd, param);

### Arguments

X, Y, XInd, YInd	Defined the same as those in function ' <i>sMBPLS_comodule</i> '.
param	A structure variable including six components:
param.maxiter, param.tol	Same as those in function ' <i>sMBPLS_comodule</i> '.
param.thrXc, param.thrYc	Two positive integers which are two thresholds for selecting features of $X$ and $Y$ in the identified md-modules.
param.thrXr, param.thrYr	$param.thrXr = param.thrYr$ . A positive integer which is a threshold for selecting samples in the identified md-modules.

### Output

success	An indicator to show if this algorithm runs successfully or not (1 for success, 0 for failure).
w, q	Weight vectors for $X, Y$ .
b, a	Super weight vectors for $X, Y$ .
T	A ( $p \times N$ ) matrix, $T(:, i)$ is the score vector for $X_i$ ( $i = 1, \dots, N$ ).
U	A ( $p \times M$ ) matrix, $U(:, i)$ is the score vector for $Y_i$ ( $i = 1, \dots, M$ ). Generally, $M = 1$ .
t, u	Super score vectors for $X, Y$ .
XX, YY	The new matrices after removing the signals of current md-modules from $X$ and $Y$ .
iter	The number of iterations to indicate when the algorithm stops.
TerminalObj	A ( $iter \times (2 + N + M)$ ) matrix in which each row records the values of all the terms in the objective function and the sum of them in each iteration.

---

---

**Description**

This function is used to incorporate all the parameters into a structure variable ‘newparams’.

**Usage**

```
newparams = sMBPLS_params(params);
```

**Arguments**

params	A structure variable including params.maxiter, params.tol, params.nfold, params.nfactor They are the same as those in the function ‘sMBPLS_comodule’.
params.thrXYr_list	This is a column vector recording the thresholds for selection samples.
params.thrXc_list	A $(1 \times N)$ cell array. <i>params.thrXc_list</i> { <i>i</i> } is a column vector, which includes several thresholds for selection features in data matrix $X_i$ ( $i = 1, \dots, N$ ).
params.thrYc_list	A $(1 \times M)$ cell array. <i>params.thrYc_list</i> { <i>i</i> } is a column vector, which includes several thresholds for selection features in $Y_i$ ( $i = 1, \dots, M$ ).

**Output**

newparams	A structure variable which is the same as the input argument ‘params’ in ‘sMBPLS_comodule’.
-----------	---

---

sMBPLS\_select\_param

Select a group of proper parameters used in sMBPLS algorithm.

---

**Description**

Using a cross-validation (CV) procedure to select a group of proper parameters.

**Usage**

```
[param_idx, cv_scores] = sMBPLS_select_param(X, Y, XInd, YInd, params);
```

**Arguments**

X, Y, XInd, YInd	Defined as those in function ‘sMBPLS_comodule’.
params	A structure variable including <i>params.param</i> , <i>params.randRowPartitions</i> , which are both the same as the output argument ‘params’ in function ‘sMBPLS_comodule’.

**Output**

param_idx	The index of selected optimal parameters in <i>params.param</i> .
cv_scores	A column vector recording the CV scores for all groups of parameters.

---

sMBPLS\_getCVscore

Calculate the CV score for one group of parameters.

---

**Description**

By using a n-fold cross-validation (CV) procedure for one group of parameters, it will obtain the corresponding CV score to assess this group of parameters. The smaller the better.

**Usage**

```
cv_score = sMBPLS_getCVscore(X, Y, XInd, YInd, param, randRowPartitions);
```

### Arguments

X, Y, XInd, YInd, param Defined the same as those in function 'sMBPLS\_algorithm'.  
randRowPartitions Defined the same as 'params.randRowPartitions' in the output of function 'sMBPLS\_comodule'.

### Output

cv\_score A score for this group of parameters 'param' in the input arguments.

---

variable\_sparse *Make the vector sparse.*

---

### Description

Make the input vector sparse based on the input threshold.

### Usage

```
[sw, err] = variable_sparse(w, thrd, msg);
```

### Arguments

w A vector.  
thrd A threshold for the degree of sparsity of  $w$ .  
msg A string representing the name of variable  $w$ .

### Output

sw The sparse vector.  
err An indicator (1 indicates it is done successfully, and 0 for not).

---

newMatrix *Obtain a new matrix.*

---

### Description

Remove the signals of the current identified md-module from the current data matrix  $X$ , and obtain a new matrix  $XX$ .

### Usage

```
XX = newMatrix(X, t, thrd, msg1, msg2);
```

### Arguments

X A matrix.  
t The latent variable for data  $X$ .  
thrd A threshold.  
msg1 A string representing the name of variable  $t$ .  
msg2 A string representing the name of variable  $X$ .

### Output

XX A new matrix removed signals from matrix  $X$ .

## 3.2 Output figures

---

sMBPLS\_plot\_XY *Provide the heatmaps for the original or reordered input matrices.*

---

### Description

Draw the heatmaps of the original input matrices  $(X, Y)$  by ignoring the input argument 'vectorForRank' or the heatmaps of reordered input matrices ordered based on the variable



'*vectorForRank*' where the signals of the identified md-module will be located in the four corners of this heatmap.

### Usage

sMBPLS\_plot\_XY(X, Y, XInd, YInd, fig, figure\_title, colormap\_type, vectorForRank);

### Arguments

X, Y, XInd, YInd      Defined as those in the function '*sMBPLS\_comodule*'.  
 fig, figure\_title, colormap\_type  
                                  Defined as those in the function '*jNMF\_plot\_X*' of jNMF.  
 vectorForRank        A structure variable including  
   vectorForRank.t     The scores used for ordering samples. It is one column of matrix  $sT$   
                                  in the output of function '*sMBPLS\_comodule*'.  
   vectorForRank.w, vectorForRank.q  
                                  The weight vectors for  $X, Y$  related to a specific md-module,  
                                  respectively. They are respectively one column of matrix  $W, Q$  in the  
                                  output of function '*sMBPLS\_comodule*'.

### Output

The heatmaps for the original or reordered input data matrices.

---

sMBPLS_plot_results	<i>Demonstrate the heatmaps of a specific md-module and the scatterplots for the correlation between the selected features.</i>
---------------------	---

---

### Description

Show the heatmaps of a selected md-module (circled in yellow lines), and the scatterplots for the correlations between the selected features.

### Usage

sMBPLS\_plot\_results(X, Y, XInd, YInd, FeatureType, fig, figure\_title, colormap\_type, vectorForRank);

### Arguments

X, Y, XInd, YInd, fig, figure\_title, colormap\_type  
                                  Defined the same as those input arguments in '*sMBPLS\_plot\_XY*'.  
 vectorForRank        A structure variable including six components. Except for the three  
                                  ones in the input argument *vectorForRank* of '*sMBPLS\_plot\_XY*',  
                                  there are also  
   vectorForRank.T     A  $(p \times N)$  matrix.  $T(:, i)$  is the score vector for  $X_i$  ( $i = 1, \dots, N$ ).  
   vectorForRank.U     A  $(p \times M)$  matrix.  $U(:, i)$  is the score vector for  $Y_i$  ( $i = 1, \dots, M$ ).  
   vectorForRank.comodule  
                                  A  $(1 \times (N + M + 1))$  cell array. *vectorForRank.comodule*{ $i$ } records  
                                  the  $i$ th feature indexes of a selected md-module. The first column is  
                                  for selected samples.  
 FeatureType         A  $(1 \times (N + M))$  cell array. *FeatureType*{ $i$ } records the name of  $i$ th  
                                  type of features (e.g., *FeatureType* = {'microRNA expression', 'CNV',  
                                  'DNA Methylation', 'Gene expression'}).

### Output

The heatmaps for certain selected md-module to demonstrate the patterns of this md-module as shown in Figure 4.

The scatterplots for the correlation between the selected features as shown in Figure 5.

---

sMBPLS_plot_distribution	<i>Demonstrate the module size distributions.</i>
--------------------------	---

---

## Description

This function provides histograms for the size distributions of all the components in the identified md-modules.

## Usage

```
sMBPLS_plot_distribution(nSample, Ind, Comodule, FeatureType, fig, figure_title);
```

## Arguments

Ind  $Ind = [XInd; YInd]$ .  $XInd$ ,  $YInd$  are the same as those in function ‘sMBPLS\_comodule’.

Comodule It is the same as that in the output of function ‘sMBPLS\_comodule’.

nSample, FeatureType, fig, figure\_title Similar with those in function ‘jNMF\_plot\_distribution’.

## Output

Histograms for the size distributions of all the features in the identified md-modules.

### 3.3 Output text files

---

Index2LabelForModuleContent	<i>Output the identified md-modules into text files.</i>
-----------------------------	--

---

## Description

This function is the same as that described in jNMF.

---

OutputModule2TXT	<i>Output the feature indexes of the identified md-modules.</i>
------------------	---

---

## Description

This function is the same as that described in jNMF.

## 4 SNPLS

SNPLS (Sparse Network-regularized PLS) is designed for one input matrix ( $X \in \mathbb{R}^{p \times n}$ ) and one response matrix ( $Y \in \mathbb{R}^{p \times m}$ ). It introduces network-regularized constraints, expressed as adjacency matrices  $A \in \mathbb{R}^{n \times n}$  of a given interaction network  $G_1$  for the features in  $X$  and/or  $B \in \mathbb{R}^{m \times m}$  of another interaction network  $G_2$  for  $Y$ . This problem is defined as,

$$\max_{g,d} cov(Xg, Yd) - \lambda_1 g^T L_X g - \lambda_2 d^T L_Y d - \lambda_3 \|g\|_1 - \lambda_4 \|d\|_1$$
$$s.t. \quad g^T g = 1, \quad d^T d = 1.$$

where  $u = Xg, v = Yd, L_X, L_Y$  are the symmetric Laplacian matrices of network  $G_1, G_2$ , respectively.

### 4.1 Algorithm

---

Run_SNPLS	<i>The main function for SNPLS.</i>
-----------	-------------------------------------

---

## Description

This is the main function for SNPLS, which integrates all the related functions to achieve it.



V	A ( $p \times nfactor$ ) matrix, $V(:,i)$ is the super score for $Y$ in the $i$ th md-module ( $i = 1, \dots, nfactor$ ).
params	Compared to 'params' as input, there are new components including <i>params.iterNumList</i> , <i>params.records</i> , <i>params.randRowPartitions</i> , <i>params.cv_scores</i> , <i>params.paramsidx_used</i> defined the same as the output argument in the function 'sMBPLS_comodule'.

SNPLS\_algorithm                      *SNPLS algorithm.*

### Description

This implements the SNPLS algorithm.

### Usage

[success, g, d, u, v, XX, YY, TerminalObj, iter] = SNPLS\_algorithm(X, Y, LX, LY, param);

### Arguments

X, Y	The same as those in function 'SNPLS_comodule'.
LX, LY	Laplacian matrices of size ( $n \times n$ ) and ( $m \times m$ ) about the interaction network for the features in data $X$ and $Y$ , respectively.
param	A structure variable including
params.tol, params.maxiter	Defined the same as those in the function 'SNPLS_comodule'.
param.thrXc, param.thrYc	Two non-negative numbers which are respectively related to the parameters $\lambda_3, \lambda_4$ in the objective function.
param.thrXNet, param.thrYNet	Two non-negative numbers. <i>param.thrXNet</i> is the parameter $\lambda_1$ in the objective function used in the network-regularized constraint about $X$ . <i>param.thrYNet</i> is similar with <i>param.thrXNet</i> used in the network about $Y$ if available.

If *param.thrXNet* = *param.thrYNet* = 0, this algorithm reduces to sMBPLS for pairwise case. If *param.thrXc* = *param.thrYc* = *param.thrXNet* = *param.thrYNet* = 0, this algorithm reduces to the standard PLS.

### Output

success, TerminalObj, iter, XX, YY	They are the same as those in the function 'sMBPLS_algorithm'.
g, d	Weight vectors for $X$ and $Y$ .
u, v	Score vectors for $X$ and $Y$ .

SNPLS\_params                      *Integrate all the parameters into an unified framework.*

### Description

This function is used to incorporate all the parameters into a structure variable 'newparams'.

### Usage

newparams = SNPLS\_params(params);

### Arguments

params	A structure variable including
params.tol, params.maxiter, params.nfold, params.nfactor, params.thrd_module	Defined the same as those in the function 'SNPLS_comodule'.
params.thrXNet_list	This is a column vector recording the parameters for the network

	constraint about $X$ in the objective function.
params.thrYNet_list	This is a column vector recording the parameters for the network constraint about $Y$ in the objective function.
params.thrXc_list	A column vector, which includes several thresholds for selecting features in data matrix $X$ .
params.thrYc_list	A column vector, which includes several thresholds for selecting features in data matrix $Y$ .

### Output

newparams	A structure variable defined the same as the input argument ' <i>params</i> ' in the function ' <i>SNPLS_comodule</i> '.
-----------	--

---

SNPLS_select_param	<i>Select a group of proper parameters used in SNPLS.</i>
--------------------	---

---

### Description

Using a cross-validation (CV) procedure to select a group of proper parameters.

### Usage

[param\_idx, cv\_scores] = SNPLS\_select\_param(X, Y, LX, LY, params);

### Arguments

X, Y, LX, LY	Defined the same as those in the function ' <i>SNPLS_algorithm</i> '.
params	A structure variable including <i>params.param</i> , <i>params.randRowPartitions</i> , which are both the same as the output argument ' <i>params</i> ' in the function ' <i>SNPLS_comodule</i> '.

### Output

param_idx	The index of the selected optimal parameters in <i>params.param</i> .
cv_scores	A column vector recording the CV scores for all groups of parameters.

---

SNPLS_getCVscore	<i>Calculate the CV score for a group of parameters.</i>
------------------	--

---

### Description

By using a n-fold cross-validation (CV) procedure for one group of parameters, it will obtain the related CV score to assess these parameters. The smaller the better.

### Usage

cv\_score = SNPLS\_getCVscore(X, Y, LX, LY, param, randRowPartitions);

### Arguments

X, Y, LX, LY, param are the same with those in the function '*SNPLS\_algorithm*'.  
 randRowPartitions is the same with '*params.randRowPartitions*' in the output of function '*SNPLS\_comodule*'.

### Output

cv_score	The score for this group of parameters ' <i>param</i> ' in the input arguments.
----------	---

---

PLS	<i>The standard PLS algorithm.</i>
-----	------------------------------------

---

### Description

This function is used to produce the initial vectors for SNPLS.

### Usage

[success, g, d, u, v] = PLS (X, Y);

## Arguments

X, Y Two input matrices for the PLS algorithm.

## Output

success An indicator to indicate whether PLS performs successfully or not (1 for success and 0 for failure).

g, d The weight vectors for input data  $X$  and  $Y$ .

u, v The score vectors for input data  $X$  and  $Y$ .

---

thresholding *Make the input vector sparse.*

---

## Description

Make the input vector sparse based on the input thresholds.

## Usage

```
sw = thresholding(w,thrd,msg);
```

## Arguments

w A vector.

thrd A threshold for the degree of sparsity of  $w$ .

msg A string for the name of variable  $w$ .

## Output

sw The sparse vector.

## 4.2 Output figures

---

SNPLS\_plot\_XY *Provide the heatmaps for the original or reordered input matrices.*

---

## Description

It is a similar function as 'sMBPLS\_plot\_XY' described in sMBPLS.

## Usage

```
SNPLS_plot_XY(X, Y, fig, figure_title, colormap_type, vectorForRank);
```

## Arguments

X, Y The same as those in the function 'SNPLS\_comodule'.

fig, figure title, colormap\_type

Defined the same as those in the function 'jNMF\_plot\_X' of jNMF.

vectorForRank A structure variable including

vectorForRank.u The scores used for ordering samples. It is one column of matrix  $U$  in the output of function 'SNPLS\_comodule'.

vectorForRank.g, vectorForRank.d

The weight vectors for  $X$  and  $Y$  related to one identified md-module, respectively. They are respectively one column of matrix  $G$ ,  $D$  in the outputs of function 'SNPLS\_comodule'.

## Output

The heatmaps for the original or reordered input data matrices.

---

SNPLS\_plot\_results *Demonstrate the heatmaps of certain identified md-module and the scatterplots for the correlation between the selected features.*

---

---

**Description**

It is a similar function with ‘*sMBPLS\_plot\_results*’ described in sMBPLS.

**Usage**

SNPLS\_plot\_results(X, Y, FeatureType, fig, figure\_title, colormap\_type, vectorForRank);

**Arguments**

X, Y, fig, figure\_title, colormap\_type  
The same as those of function ‘*SNPLS\_plot\_XY*’.

vectorForRank  
A structure variable containing  
vectorForRank.g, vectorForRank.d, vectorForRank.u  
Defined the same as those in the function ‘*SNPLS\_plot\_XY*’.

vectorForRank.comodule  
A (1 × 3) cell. It has the similar meaning with that in the function  
‘*sMBPLS\_plot\_results*’.

vectorForRank.v  
The score vector for data *Y*. It is one column of matrix *V* in the output  
of function ‘*SNPLS\_comodule*’.

FeatureType  
Similar definition with that in the function ‘*sMBPLS\_plot\_results*’  
with  $N = M = 1$ .

**Output**

The heatmaps for certain selected md-module to demonstrate the patterns of this md-module.  
The scatterplots for the correlation between the selected features.

---

SNPLS_plot_distribution	<i>Demonstrate the module size distributions.</i>
-------------------------	---

---

**Description**

This function provides histograms for the size distributions of all the components in the identified md-modules, which has the similar input arguments and output results.

**Usage**

SNPLS\_plot\_distribution(nSample, Ind, Comodule, FeatureType, fig, figure\_title);

### 4.3 Output text files

---

Index2LabelForModuleContent	<i>Output the identified md-modules into text files.</i>
-----------------------------	--

---

**Description**

This function is the same as that described in jNMF.

---

OutputModule2TXT	<i>Output the feature indexes of identified md-modules.</i>
------------------	---

---

**Description**

This function is the same as that described in jNMF.

## 5 Input data

To facilitate the usage, MIA package implements the four methods using the same structure variable to describe the input data. This variable, named *Input*, includes the following components:

*Input.data*: A matrix storing all the multi-dimensional data sequentially (e.g.,  $Input.data = [X_1, \dots, X_N]$ ). Each row corresponds to the genomics features of a specific sample. Each type of genomic data is assigned its own set of columns.

*Input.XBlockInd*: A matrix of size  $N \times 2$ . The two elements in  $i$ th row give the start and end column indexes in *Input.data* for the  $i$ th  $X$  matrix ( $i = 1, \dots, N$ ).

*Input.YBlockInd*: A matrix storing the response matrix  $Y$  for both sMBPLS and SNPLS. Its format is similar to *Input.XBlockInd*.

*Input.netAdj*: The symmetric adjacency matrix of a given network used for SNMNMF and SNPLS, where the features have the same order as in *Input.data*. This network combines the interactions between and within the variables in multiple types of variables. The element of this matrix equals to 1 for linked features in the network, and 0 otherwise.

*Input.SampleLabel*: A vector recording the labels of samples.

*Input.FeatureLabel*: A vector recording the feature names in *Input.data*. The  $i$ th label corresponds to the  $i$ th feature in *Input.data*.

*Input.FeatureType*: A vector recording the feature types in *Input.data*. Here we give an example:  $Input.FeatureType = \{ \text{'Gene expression'}, \text{'microRNA expression'}, \text{'DNA methylation'} \}$ .

*Input.params*: A structure variable, storing all the parameters used in MIA.

For these four methods, there are three common parameters, including

- *Input.params.NCluster*: A pre-defined number of md-modules. For example, we may set  $Input.params.NCluster = 20$ .
- *Input.params.maxiter*: The maximal iteration times in each algorithm. For example, we may set  $Input.params.maxiter = 100$ .
- *Input.params.tol*: The precision for convergence of each algorithm. For example, we may set  $Input.params.tol = 10^{-6}$ .

For jNMF, there are two specific parameters:

- *Input.params.nloop*: The number of repeating times to run this algorithm. To obtain a robust and optimal solution, this algorithm is run for multiple times repeatedly, and the solution with the minimal value of objective function is accepted. For example, we may set  $Input.params.nloop = 50$ .
- *Input.params.thrd\_module*: A non-negative vector of size  $1 \times (N + 1)$  to select features in md-modules.  $Input.params.thrd\_module(i+1)$  is the threshold for selecting the  $i$ th type of features in *Input.data* ( $i = 1, \dots, N$ ). The first one is for selecting samples. The larger they are, the smaller number of features are selected. Users can set it based on the size of md-modules they prefer to identify. For example, we may  $Input.params.thrd\_module = ones(1, N + 1)$ .

For SNMNMF, except for *Input.params.nloop* and *Input.params.thrd\_module*, there are also:

- *Input.params.thrNet11*, *Input.params.thrNet12*, *Input.params.thrNet22*: The three non-negative numbers are set for the parameters respectively related to the network constraints about network  $A_{11}$ ,  $A_{12}$ ,  $A_{22}$  in the objective function, where  $A_{11}$ ,  $A_{22}$  are respectively the adjacency matrices for the interaction networks within the features in data matrix  $X_1$ ,  $X_2$ ;  $A_{12}$  is for the interaction network between the two types of features. User can choose which networks they prefer to use in the framework by setting the corresponding parameters. For example, if  $Input.params.thrNet11 = 0$ , the network  $A_{11}$  will not be used.



- *Input.params.thrXr*, *Input.params.thrXc*: The two non-negative numbers are set for the row related (or  $W$ ), and column related ( $H_i$ ) terms respectively in the objective function. It controls the degree of sparsity of matrix  $W$ ,  $H$ . For example, we may set *Input.params.thrXr* = 10, *Input.params.thrXc* = 10.

For sMBPLS, there are:

- *Input.params.nfold*: A positive number used for  $n$ -fold cross-validation (CV) procedure. Generally, we set *Input.params.nfold* = 5 or = 10. This method applies CV procedure to select a proper group of parameters from all the combinations of these parameter lists described below.
- *Input.params.thrXYr\_list*: A column vector with positive integers. They are candidates for thresholds in order to select samples in md-modules. For example, we may set *Input.params.thrXYr\_list* = [20; 30].
- *Input.params.thrXc\_list*, *Input.params.thrYc\_list*: Two row vectors of size  $1 \times N$ ,  $1 \times M$  with positive integers to control the degree of sparsity for the weight variables of input data  $X$ ,  $Y$ , respectively. For example, we may set *Input.params.thrXc\_list* = repmat([20; 30], 1,  $N$ ), *Input.params.thrYc\_list* = repmat([20; 30], 1,  $M$ ), where  $N = \text{size}(\text{Input.XBlockInd}, 1)$ ,  $M = \text{size}(\text{Input.YBlockInd}, 1)$ .

For SNPLS, there are:

- *Input.params.nfold*: It is the same as that in sMBPLS.
- *Input.params.thrXc\_list*, *Input.params.thrYc\_list*: They have the same meaning as those in sMBPLS for the situation of  $N = 1$ ,  $M = 1$ . Thus, they are defined as column vectors. For example, we may set *Input.params.thrXc\_list* = [0.01; 0.03; 0.05], *Input.params.thrYc\_list* = [0.1; 0.3; 0.5].
- *Input.params.thrXNet\_list*, *Input.params.thrYNet\_list*: The two column vectors with non-negative values. They have the similar function with *Input.params.thrNet11* in SNMNMF. *Input.params.thrXNet\_list*, *Input.params.thrYNet\_list* are respectively for the networks within the features in input data  $X$ , and response data  $Y$ . For example, we may set *Input.params.thrXNet\_list* = [1; 5]; *Input.params.thrYNet\_list* = [1; 5].
- *Input.params.thrd\_module*: It is a non-negative matrix of size  $3 \times 2$ . The first column *Input.params.thrd\_module*( $i$ , 1) is the threshold for selecting the  $i$ th feature in *Input.data* ( $i = 1, 2, 3$ ). And the second column *Input.params.thrd\_module*( $i$ , 2) is a percentage in case of no features selected using the threshold. The first row is for selecting samples. The larger the thresholds are, the smaller number of features are selected. Users can set it based on the size of md-modules they prefer to identify. For example, we may set *Input.params.thrd\_module* = [1, 0.5; 1, 0.5; 1, 0.5].

In addition, for the components that are not used in certain methods (e.g., *Input.YBlockInd* in jNMF and SNMNMF and *Input.netAdj* in jNMF and sMBPLS), users can set them null or just ignore them.

With this data structure, MIA is able to partition *Input.data* into corresponding data matrices as input for each method automatically.

Next, we provide an example for constructing the input data used in SNMNMF. Suppose that one wants to identify 50 microRNA-gene co-modules by integrate gene expression profiles ( $X_1 \in \mathbb{R}^{385 \times 12456}$ ) and micro-RNA expression profiles ( $X_2 \in \mathbb{R}^{385 \times 559}$ ) across the same set of samples, as well as the gene interaction network  $G_1$ , gene-microRNA interaction network  $G_2$ . The network  $G_1$  can be expressed by the adjacency matrix  $A_{11} = (a_{ij})_{12456 \times 12456}$ , where  $a_{ij} = 1$  if gene  $i$  and gene  $j$  is linked in the network  $G_1$ . Similarly,  $G_2$  is expressed by the

adjacency matrix  $A_{12} \in \mathbb{R}^{12456 \times 559}$ . If the microRNA interaction network is not available, the corresponding adjacency matrix  $A_{22}$  is defined as  $A_{22} = \text{zeros}(559, 559)$ .

Then, we could define the input data *Input* as below:

---

```

Input.data = [X1, X2];
Input.XBlockInd = [1, 12456; 12457, 13015];
Input.YBlockInd = [ ];
Input.netAdj = [A11, A12; A12T, A22];
Input.SampleLabel = {'TCGA-24-1105-01A';...;'TCGA-13-0793-01A'};
Input.FeatureLabel = {'SFRS8';...;'SCN3A';'hsa-mir-488';...;'hsa-mir-874'};
Input.FeatureType = {'Gene', 'miRNA'};
Input.params.NCluster = 50;
Input.params.maxiter = 100;
Input.params.tol = 10-6;
Input.params.nloop = 5;
Input.params.thrd_module = [1,0.5;1,0.5;1,0.5];
Input.params.thrNet11 = 10-4; Input.params.thrNet12 = 0.01; Input.params.thrNet22 = 0;
Input.params.thrXr = 10; Input.params.thrXc = 10;

```

---

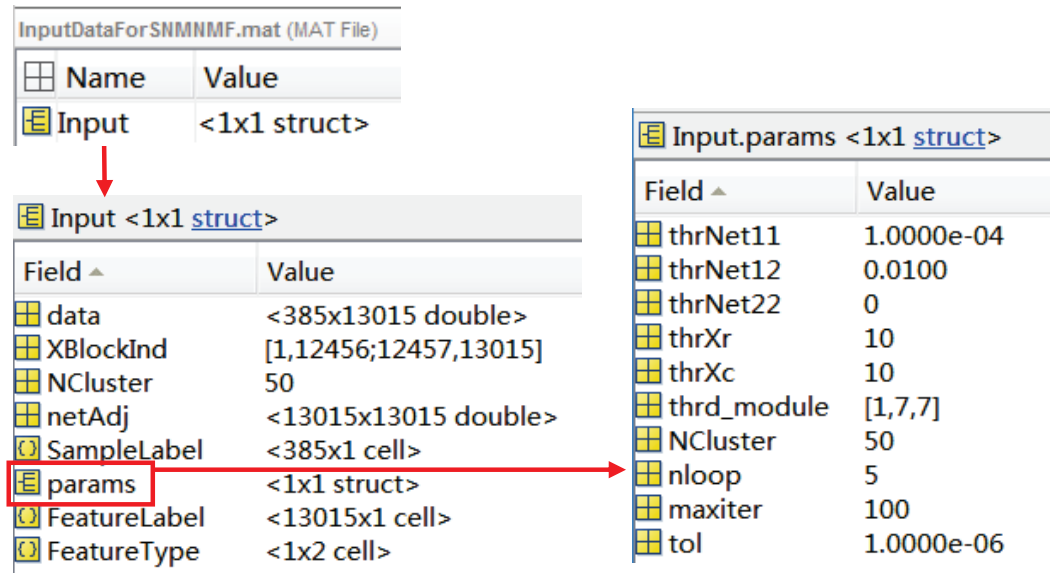


Figure 1: Illustration of an example of the input data for *SNMNMf*.

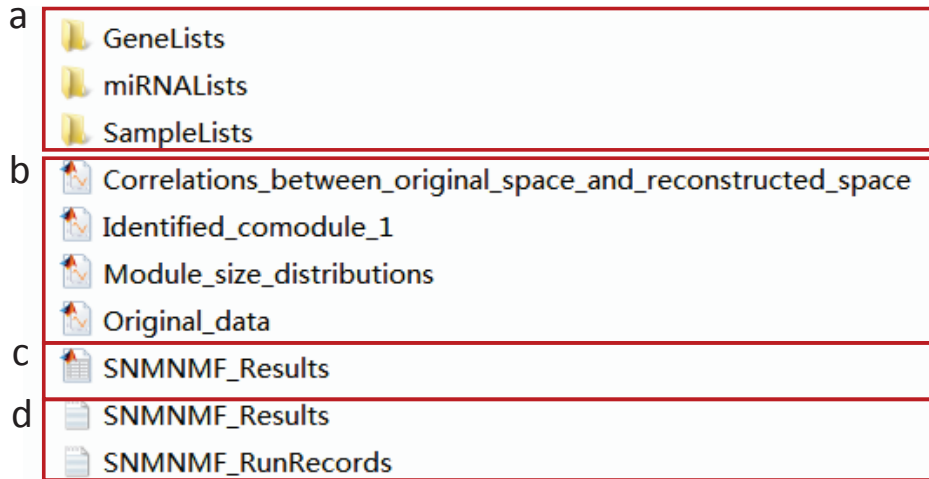


Figure 2: An overview of the output results for *SNMNMf*. The details about each part are shown in Figure 3 and Figure 4.

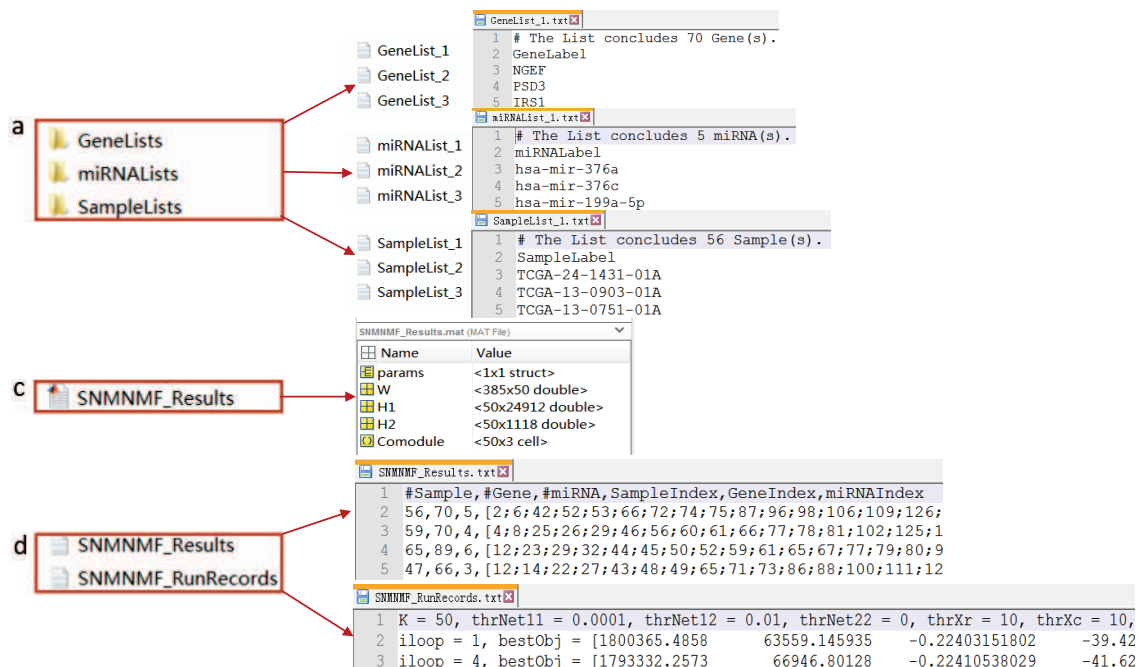


Figure 3: The details about the output files shown in Figure 2. (a) In each folder, there are a number of text files, each of which records one type of components in one identified co-modules. (c) A MATLAB data file storing the computation results, including the factorized matrices  $W$ ,  $H1$ ,  $H2$ , the 50 identified co-modules, and the parameters used in this method. (d) The first text file records the feature indexes of all the identified md-modules, in which the first three numbers are the number of samples, genes, microRNAs in one identified md-modules, and the next three columns show the indexes of selected samples, genes and microRNA, respectively. Each list are included in the square brackets. The second text file records some information during the iterations. The first line shows the parameters used in SNMNMf. The rest lines show the changes of objective function during multiple-round running. It just records the results better than the previous round. *bestObj* stores values of the terms in the objective function in the '*iloop*'th round, and *sum\_Obj* is the sum of these terms.

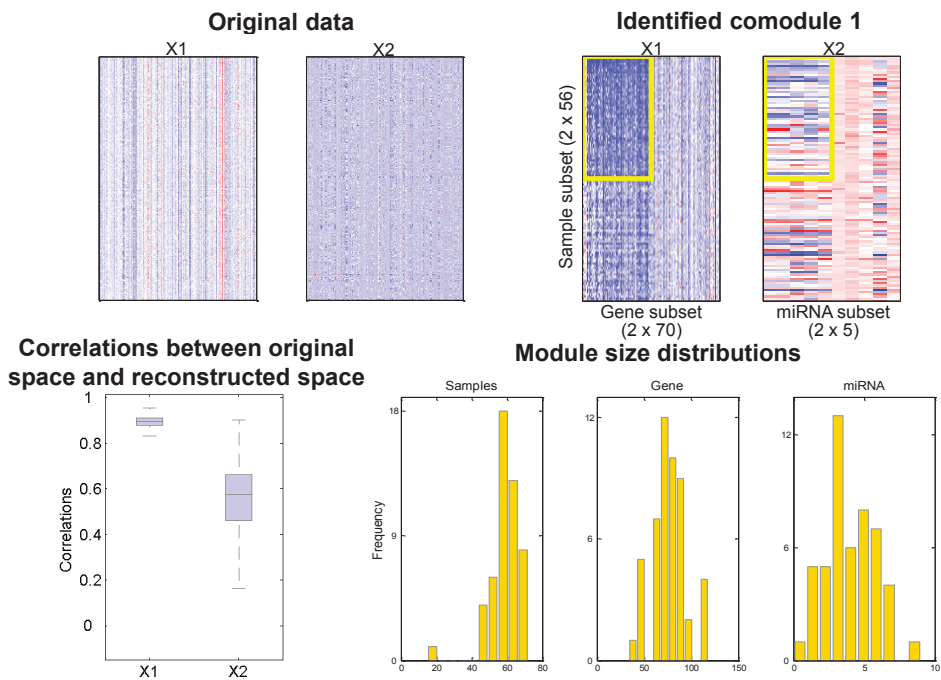


Figure 4: An example for the output figures in Figure 2b.

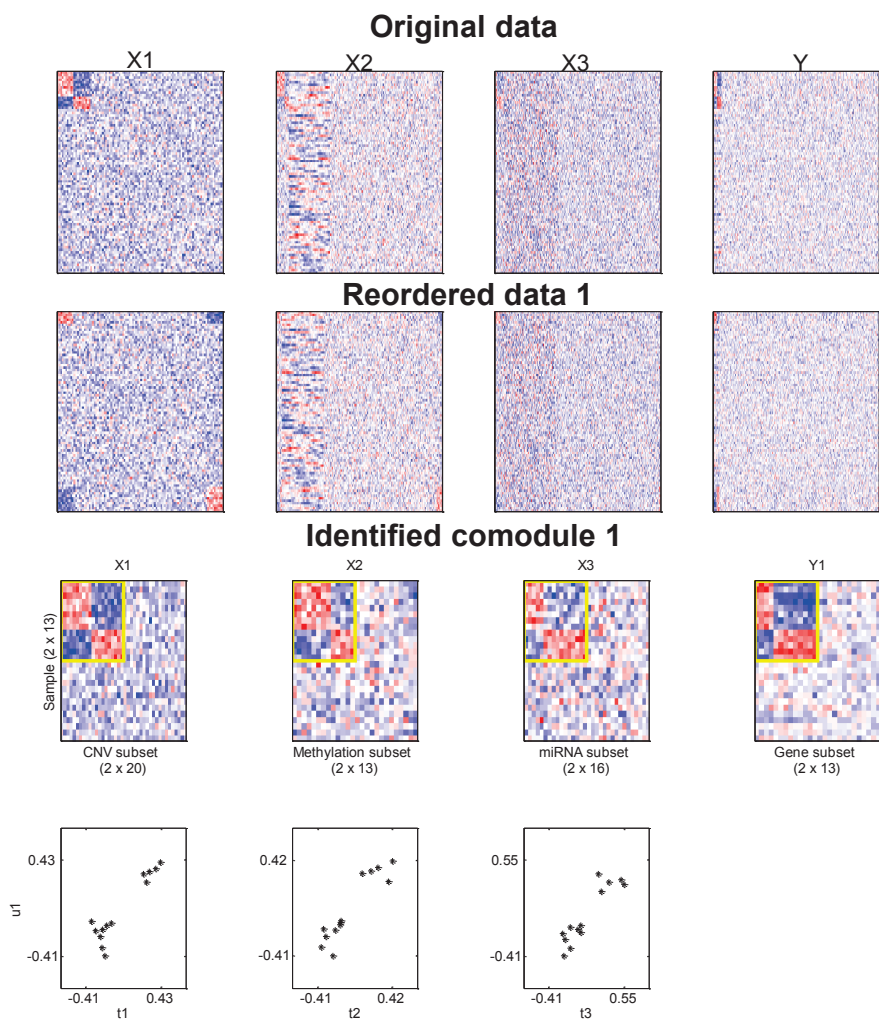


Figure 5: An example for the output figures for sMBPLS.

## 6 Guide for the MIA users without a MATLAB license

### Guide:

1. Operating System requirements: Windows 64-bit.
2. Download and install the Windows 64-bit version of the MATLAB Runtime for R2015b from the MathWorks Web site by navigating to <http://www.mathworks.com/products/compiler/mcr/index.html>.
3. Prepare the input data and store them in the path “./MIA/InputData/”. In this folder, we provide the input data for each method as examples. For each method, there are two Excel files (one is for input data and another one is for input parameters), each of which includes several sheets. Users need to arrange their data in the same way as those example Excel files in this folder we provided. Note that, each sheet is renamed as the corresponding variable name as described in our manuscript.
4. Open the Command Prompt (cmd.exe). Set the current path as where the MIA package is located, e.g., “D:/MIA”.

5. Produce MATLAB data files (\*.mat) by running PreInputData.exe. Type the command as below:

```
D:\MIA> PreInputData.exe ./InputData/DataForjNMF.xlsx ./InputData/ParametersForjNMF.xlsx
./InputData/InputDataForjNMF.mat jNMF
or
D:\MIA> PreInputData.exe ./InputData/DataForSNMNMF.xlsx ./InputData/ParametersForSNMNMF.xlsx
./InputData/InputDataForSNMNMF.mat SNMNMF
or
D:\MIA> PreInputData.exe ./InputData/DataForsMBPLS.xlsx ./InputData/ParametersForsMBPLS.xlsx
./InputData/InputDataForsMBPLS.mat sMBPLS
or
D:\MIA> PreInputData.exe ./InputData/DataForSNPLS.xlsx ./InputData/ParametersForSNPLS.xlsx
./InputData/InputDataForSNPLS.mat SNPLS
```

The first two parameters are the file names storing input data and parameters; the next parameter is the output file name. The last one is the selected method. The produced new data files (e.g., “InputDataForjNMF.mat”) are saved in the path “./MIA/InputData/”.

6. Run MIA.exe. Type the command as below:

```
D:\MIA> MIA.exe ./InputData/InputDataForjNMF.mat jNMF
or
D:\MIA> MIA.exe ./InputData/InputDataForSNMNMF.mat SNMNMF
or
D:\MIA> MIA.exe ./InputData/InputDataForsMBPLS.mat sMBPLS
or
D:\MIA> MIA.exe ./InputData/InputDataForSNPLS.mat SNPLS
```

The first one is about the input data file and the second one is about the selected method. For each method, the results are saved their own directory. For example, the results of running jNMF are saved in “./MIA/jNMF/jNMF\_Results/”.