# User guide for the MIIC algorithm command line version

Nadir Sella,<sup>1,2</sup> Louis Verny,<sup>1,2,3</sup> Séverine Affeldt,<sup>1,2,4</sup> Hervé Isambert,<sup>1,2,†</sup>

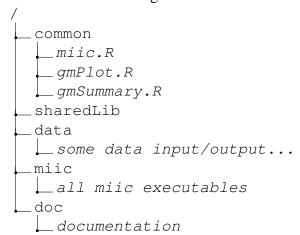
<sup>1</sup>Institut Curie, PSL Research University, CNRS, UMR168, 26 rue dUlm, 75005 Paris, France <sup>2</sup>Sorbonne Universités, UPMC Univ Paris 06, 4, Place Jussieu, 75005 Paris, France

<sup>3</sup> Current address: CENOVA SAS, 100 Avenue Charles De Gaulle, 92200 Neuilly Sur Seine <sup>4</sup> Current address: ICAN, UPMC, Paris, France

<sup>†</sup>To whom correspondence should be addressed; E-mail: herve.isambert@curie.fr

The main folder contains the scripts and source code for the reconstruction of networks starting from observation data. The structure is the following:

The directories are organized as follow:



## Package requirements

To launch the miic.R script you need to have R installed on your machine, along with some packages that are available in the CRAN repository.

Rpackages MASS, getopt, plotrix, methods, igraph, ppcor, bnlearn

### Calling the inference methods with miic

You can call the inference methods through the *miic*.*R* script.

#### Overview

**main** ∼/common/miic.R

**lib** ∼/common/lib/...

Arguments (mandatory: \*)

- -i \* file path of the input dataset<sup>1</sup>
- -o \* directory path for the output of the inference method<sup>2</sup>
- **-d** steps to perform ('1,2,3,4' or '1,2' or '1,3' *etc...*) default: '1,2,3,4'<sup>3</sup>
- -p parameters for the inference method (see the following subsections). The value expected here is of type character: 'param<sub>1</sub>:value<sub>1</sub>,param<sub>2</sub>:value<sub>2</sub> etc...'
- -t file path to the true edges; used during the summary step<sup>4</sup>
- -s file path to the category order used during the *summary* step<sup>5</sup>
- -x file path to the blackbox file containing the edges to be removed at the beginning of the reconstruction step.<sup>6</sup>

<sup>&</sup>lt;sup>1</sup>The input dataset should be a tab separated table, with column names but no row names. Missing values should be indicated with NA. Each column corresponds to a categorized variable and each row to one sample.

<sup>&</sup>lt;sup>2</sup>To prevent from overwriting existing results, if the output directory already exists, the skeleton inference step returns a message and stops.

<sup>&</sup>lt;sup>3</sup>(1) skeleton, (2) orientation, (3) summary, (4) plot

<sup>&</sup>lt;sup>4</sup>The true edges file has two space-separated columns. Each line corresponds to one edge. The orientation is  $col1 \rightarrow col2$ .

 $<sup>^5</sup>$ The orientation is  $col1 \rightarrow col2$ . This files provides information about how to consider the different states of categorical variables. It will be used to compute the signs of the edges (using spearman correlation coefficient) by ranking the levels of each variable according to the order given in the file. This file is necessary (except for numerical variables) to obtain edge colors corresponding to the signs of their partial correlations (positive in red, negative in blue). If it is not possible or desirable to order the states of some variables, the column "levels\_increasing\_order" can be left empty for these variables. The edges involving those variables are then colored in gray in the reconstructed network. (NB: in this case, the field separator is still needed between the node name and the empty "levels\_increasing\_order" cell in the category order file).

<sup>&</sup>lt;sup>6</sup>It must be formatted as a two-column file, Node1 Node2, with a field separator between them.

- -I file path to the layout of each vertex; used during the *plot* step<sup>7</sup>
- -c if given, edges will be filtered according to their confidence ratio. It needs two parameters, described in . To use a different threshold without shuffling the data a second time, an identical command line with a different confidence ratio threshold can be used. The program will keep the edges satisfying the new confidence ratio, using the previously calculated mutual information
- -v if given, detailed information is given in the log file and command shell
- -n if given, the graphml file is created in the output folder<sup>8</sup>

#### An example:

Rscript miic.R -i ../data/asiaDataset.txt -o ../data/asiaNetwork -m miic -p cpx:nml,efn:-1,lat:yes,prg:yes -c csh:100.ccr:0.01

When calling the available inference methods with miic.R, the 'p' option can be used to indicate the chosen parameters. The value expected for this option is of type character: ' $param_1:value_1,param_2:value_2$  etc...'. The possible  $param_i$  and  $value_i$  for each method are detailed in the following subsections.

#### Option '-p' for miic

```
cpx formula used to compute the complexity term ['mdl'<sup>9</sup> or 'nml'<sup>10</sup>] default: nml (Ex.: -p '...,cpx:mdl,...')
```

**lat** should the network be reconstructed under the hypothesis that some variables might not be observed? ['yes' or 'no'] default: no (*Ex.: -p ...,lat:yes,...*)

**prg** should the network be oriented using the propagation rule? ['yes' or 'no'] default: yes (*Ex.: -p ...,prg:yes,...*)

efn number of uncorrelated samples

default: number of rows of the input dataset (Ex.: -p ...,efn:1000,...)

A '-p' example: -p cpx:mdl,efn:1000

 $<sup>^{7}</sup>$ The layout file has three tab separated columns, the first column being optional. Each line corresponds to the (x,y) coordinates of each vertex. The first column can contain the label of the vertex as indicated in the colnames of the input dataset table. The order in which the coordinates are given also corresponds to the order of the colnames of the input dataset table.

<sup>&</sup>lt;sup>8</sup>If the layout file is provided, the network will be stored in a xgmml file format, that allows the node positioning in the Cytoscape tool. We have noticed that in this case edges are note coloured, we will soon try to solve it.

<sup>&</sup>lt;sup>9</sup>Maximum description length

<sup>&</sup>lt;sup>10</sup>Normalized Maximum Likelihood

#### Option '-c' for miic

**csh** number of random shuffling of the input dataset, in order to get the random mutual information between miic inferred edges (*Ex.: -c ...,csh:10,...*)

ccr confidence ratio used as a threshold for filtering the edges. (Ex.: -c ...,ccr:0.01,...)

A '-c' example: -c csh:100,ccr:0.01

#### Viewing inferred networks

The inferred networks can be viewed in pdf format, automatically generated with igraph (http://igraph.org/) or manipulated for better display using cytoscape (http://www.cytoscape.org). The files are located in the following directories:

- Unfiltered network: the pdf and graphml files are located in the output directory set by the -o entry in the command line
- Filtered networks (using -c option with miic): the pdf and graphml file are located in the subdirectory 'shuffle\_[cshValue]/filtered\_network\_[ccrValue]', which can be found in the output directory set by the -o entry in the command line

In order to visualize the network in a correct, pleasant and interactive way, we recommend the utilization of Cytoscape tool, version 3.1.0 or later (http://www.cytoscape.org/). Cytoscape is available for Windows, Linux and OsX.

In order to load the network you have to go through the following steps:

- 1. Import the network: File⇒Import⇒Network⇒File, and select the graphml file
- 2. Import the style: File⇒Import⇒Styles, and select the miic\_style.xml file present in this folder
- 3. Select the loaded style: under the "Style" panel present in "Control Panel" select miic\_style