

The mlica Package

February 16, 2008

Title Independent Component Analysis using Maximum Likelihood

Version 0.6

Author Andrew Teschendorff

Description An R code implementation of the maximum likelihood (fixed point) algorithm of Hyvaerinen, Karhuna and Oja for independent component analysis.

Maintainer Andrew Teschendorff <aet21@hutchison-mrc.cam.ac.ac>

License GPL version 2.0

R topics documented:

CheckStability	1
PriorNormPCA	3
SortModes	4
mlica	5
mlicaMAIN	6
proposeNCP	7
simMAdata	8
Index	10

CheckStability *Tests stability of inferred ICA modes.*

Description

Performs a correlation test to see which of the inferred ICA modes are reproducible across multiple runs using different random initialisations. Returns a set of consensus ICA modes and stability scores for each following the algorithm of Chiappetta,...et.al (2004).

Usage

```
CheckStability(a.best.l, corr.th)
```

Arguments

<code>a.best.l</code>	List of <code>a.best</code> objects from <code>mlica</code> runs.
<code>corr.th</code>	Correlation threshold to use to decide whether a mode is reproducible.

Value

A list with the following components

<code>consS</code>	Consensus source matrix with columns labeling the consensus ICA modes. Has same number of rows as <code>a.best\$S</code> .
<code>consA</code>	Consensus mixing matrix with rows labeling the consensus ICA modes.
<code>stabM</code>	Vector of same length as <code>consM</code> giving the stability measures of each consensus ICA mode. Stability or reproducibility measures are given as fractions, that is, the number of times the ICA mode correlates with one of the other runs at threshold level <code>corr.th</code> divided by the number of runs (length of <code>a.best.l</code>).

Author(s)

Andrew Teschendorff <aet21@cam.ac.uk>

References

- 1 Hyvaerinen A., Karhunen J., and Oja E.: *Independent Component Analysis*, John Wiley and Sons, New York, (2001).
- 2 Kreil D. and MacKay D. (2003): *Reproducibility Assessment of Independent Component Analysis of Expression Ratios from DNA microarrays*, Comparative and Functional Genomics **4** (3),300–317.
- 3 Liebermeister W. (2002): *Linear Modes of gene expression determined by independent component analysis*, Bioinformatics **18**, no.1, 51–60.
- 4 Chiappetta P., Roubaud MC. and Torresani B.: *Blind source separation and the analysis of microarray data*, J. Comput. Biol. 2004; 11(6):1090–109.

Examples

Description

This function performs a simple PCA analysis to aid in threshold setting and noise removal.

Usage

```
PriorNormPCA(X)
```

Arguments

X	Data Matrix (need not be normalised). Subsequent ICA seeks independent modes as independent distributions with values "down the rows".
---	--

Details

This function performs a simple PCA analysis and is used prior to application of the main ICA algorithm. The objective of the prior PCA is to help determine the dimensionality of a subspace on which the further ICA converges. The convention used here is that the rows of X label the space over which independent components are sought. For a typical microarray application in which ICA is being used as a generative model for gene expression, rows should label genes and columns should label samples. If, however, ICA is to be used as an unsupervised projection pursuit algorithm, rows should label samples and columns genes. For the latter application, the number of genes should be less than the number of samples.

Value

A list with following components:

X	Normalised data matrix with the mean of each column set to zero.
Dx	Eigenvalues in a diagonal matrix.
Ex	Eigenvectors

Author(s)

Andrew Teschendorff <aet21@cam.ac.uk>

References

- 1 Hyvaerinen A., Karhunen J., and Oja E.: *Independent Component Analysis*, John Wiley and Sons, New York, (2001).
- 2 Kreil D. and MacKay D. (2003): *Reproducibility Assessment of Independent Component Analysis of Expression Ratios from DNA microarrays*, *Comparative and Functional Genomics* **4** (3),300–317.

- 3 Liebermeister W. (2002): *Linear Modes of gene expression determined by independent component analysis*, *Bioinformatics* **18**, no.1, 51–60.

Examples

SortModes *Sorting of ICA modes*

Description

Sorts inferred ICA modes using two criteria: Relative data power or the Liebermeister criterion, which is based on a measure that is a weighted linear combination of non-gaussianity and data variance measures.

Usage

```
SortModes(a.best, c.val = 0.25)
```

Arguments

<code>a.best</code>	The output object of <code>mlica</code> .
<code>c.val</code>	A parameter to control the relative weight of the two measures when using the Liebermeister criterion. Should be between 0 (pure data variance measure) and 1 (pure non-gaussianity).

Value

A list with components:

<code>a.best</code>	The output of <code>mlica</code> .
<code>rdp</code>	The relative data power values obtained for each independent component.
<code>lbm</code>	The Liebermeister contrast value for each component.

Author(s)

Andrew Teschendorff <aet21@cam.ac.uk>

References

- 1 Hyvaerinen A., Karhunen J., and Oja E.: *Independent Component Analysis*, John Wiley and Sons, New York, (2001).
- 2 Kreil D. and MacKay D. (2003): *Reproducibility Assessment of Independent Component Analysis of Expression Ratios from DNA microarrays*, *Comparative and Functional Genomics* **4** (3),300–317.
- 3 Liebermeister W. (2002): *Linear Modes of gene expression determined by independent component analysis*, *Bioinformatics* **18**, no.1, 51–60.

Examples

mlica	<i>Maximum likelihood implementation of Independent Component Analysis</i>
-------	--

Description

This function performs ICA using a maximum likelihood framework and takes as arguments parameters to control the number of algorithm runs and convergence criteria.

Usage

```
mlica(prNCP, nruns = 10, tol = 1e-04, maxit = 300, fail.th = 5, learn.mu = 1)
```

Arguments

prNCP	The output object from <code>proposeNCP</code> .
nruns	The number of converged algorithm runs sought (function returns the best solution according to the log-likelihood value).
tol	Tolerance level for establishing convergence of run.
maxit	Maximum number of iterations to allow per run.
fail.th	A threshold on the number of consecutive runs that fail to converge.
learn.mu	Learning parameter for fixed point algorithm (note that this need not be changed since it has already been optimised).

Value

A list with following components:

A	Estimate of the mixing matrix.
B	Estimate of the inverse mixing matrix.
S	Estimate of the source matrix.
X	Normalised data matrix.
ncp	Number of independent components.
NC	Binary number specifying whether best run converged or not.(=1 indicates convergence,=0 indicates no convergence).
LL	Log likelihood value of best run.

Author(s)

Andrew Teschendorff <aet21@cam.ac.uk>

References

- 1 Hyvaerinen A., Karhunen J., and Oja E.: *Independent Component Analysis*, John Wiley and Sons, New York, (2001).
- 2 Kreil D. and MacKay D. (2003): *Reproducibility Assessment of Independent Component Analysis of Expression Ratios from DNA microarrays*, *Comparative and Functional Genomics* **4** (3),300–317.
- 3 Liebermeister W. (2002): *Linear Modes of gene expression determined by independent component analysis*, *Bioinformatics* **18**, no.1, 51–60.
- 4 Chiappetta P., Roubaud MC. and Torresani B.: *Blind source separation and the analysis of microarray data*, *J. Comput. Biol.* 2004; 11(6):1090–109.

Examples

```
data(simMadata);
dataX <- simMadata[[1]];
prPCA <- PriorNormPCA(dataX);
prNCP <- proposeNCP(prPCA,0.1);
a.best.l <- list();
for( i in 1:5){
  a.best.l[[i]] <- mlica(prNCP,nruns=5);
}
checkICA <- CheckStability(a.best.l,0.7);
sourceS <- simMadata[[3]];
print(cor(a.best.l[[1]]$S,sourceS));
sModes <- SortModes(a.best.l[[1]],c.val=0.5);
```

mlicaMAIN

Main engine function that implements the fixed point algorithm for maximum likelihood inference of ICA modes.

Description

See references for detailed description.

Usage

```
mlicaMAIN(prNCP, tol = 1e-04, maxit = 300, mu = 1)
```

Arguments

prNCP	The output object of proposeNCP.
tol	Tolerance level for convergence.
maxit	Maximum number of iterations to allow for convergence.
mu	Learning paramter for fixed point algorithm. This has already been optimised.

Value

A list with following components:

A	Estimate of the mixing matrix.
B	Estimate of the inverse mixing matrix.
S	Estimate of the source matrix.
X	Normalised data matrix.
nCP	Number of independent components.
NC	Binary number specifying whether best run converged,0, or not,1.
LL	Log likelihood value of best run.

Author(s)

Andrew Teschendorff <aet21@cam.ac.uk>

References

- 1 Hyvaerinen A., Karhunen J., and Oja E.: *Independent Component Analysis*, John Wiley and Sons, New York, (2001).
- 2 Kreil D. and MacKay D. (2003): *Reproducibility Assessment of Independent Component Analysis of Expression Ratios from DNA microarrays*, *Comparative and Functional Genomics* **4** (3),300–317.
- 3 Liebermeister W. (2002): *Linear Modes of gene expression determined by independent component analysis*, *Bioinformatics* **18**, no.1, 51–60.

Examples

proposeNCP *Number of independent components proposal function*

Description

This function takes the output of `PriorNormPCA` and returns for a given threshold the number of components to be inferred for subsequent ICA.

Usage

```
proposeNCP(prPCA, thresh = 0.1)
```

Arguments

prPCA	The output object from <code>PriorNormPCA</code> .
thresh	Threshold on eigenvalues.

Value

A list with following components:

X	Normalised data matrix.
x	Normalised data matrix projected onto selected subspace.
pEx	Selected eigenvectors defining subspace for projection.
pCorr	Projected correlation matrix.
ncp	The dimension of the selected subspace(=number of independent components to be inferred with subsequent ICA).

Author(s)

Andrew Teschendorff (aet21@cam.ac.uk)

References

- 1 Hyvaerinen A., Karhunen J., and Oja E.: *Independent Component Analysis*, John Wiley and Sons, New York, (2001).
- 2 Kreil D. and MacKay D. (2003): *Reproducibility Assessment of Independent Component Analysis of Expression Ratios from DNA microarrays*, *Comparative and Functional Genomics* **4** (3),300–317.
- 3 Liebermeister W. (2002): *Linear Modes of gene expression determined by independent component analysis*, *Bioinformatics* **18**, no.1, 51–60.

Examples

simMAdata	<i>Simulated Microarray data for testing purposes</i>
-----------	---

Description

This data set contains a mock microarray data set of 1000 genes and 60 samples where the data has been generated from an underlying Independent Component Analysis model of 5 supergaussian modes.

Usage

```
data(simMAdata)
```

Format

A list with three components. First, second and third components are the data matrix, mixing matrix and source matrix, respectively.

References

Hyvaerinen A., Karhunen J. and Oja E. (2001) *Independent Component Analysis*. New York: Wiley.

Index

*Topic **cluster**

CheckStability, 1

mlica, 5

PriorNormPCA, 2

proposeNCP, 7

SortModes, 4

*Topic **datasets**

simMAdata, 8

*Topic **internal**

mlicaMAIN, 6

CheckStability, 1

mlica, 5

mlicaMAIN, 6

PriorNormPCA, 2

proposeNCP, 7

simMAdata, 8

SortModes, 4