

## **FABIA: Factor Analysis for Biclust Acquisition** **— *Manual for the R package* —**

**Sepp Hochreiter**

Institute of Bioinformatics, Johannes Kepler University Linz  
Altenberger Str. 69, 4040 Linz, Austria  
*hochreit@bioinf.jku.at*

**Version 0.1.3, April 23, 2010**

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Getting Started: FABIA</b>	<b>3</b>
2.1	Quick start : without outputs and plots . . . . .	3
2.2	Test on Toy Data Set . . . . .	5
2.3	Demos . . . . .	6
<b>3</b>	<b>The FABIA Model</b>	<b>7</b>
3.1	Model Assumptions . . . . .	7
3.2	Sparse Coding and Laplace Prior . . . . .	9
3.3	Model Selection . . . . .	10
3.3.1	Variational Approach for Sparse Factors . . . . .	10
3.3.2	New Update Rules for Sparse Loadings . . . . .	10
3.3.3	Extremely Sparse Priors . . . . .	11
3.4	Data Preprocessing and Initialization . . . . .	11
3.5	Information Content of Biclusters . . . . .	12
3.6	Extracting Members of Biclusters . . . . .	13
3.7	C implementation of FABIA . . . . .	13
<b>A</b>	<b>Class Factorization</b>	<b>13</b>
<b>B</b>	<b>Biclustering and Matrix Factorization Methods</b>	<b>15</b>
B.1	fabi . . . . .	15
B.2	fabia . . . . .	19
B.3	fabiap . . . . .	25
B.4	fabias . . . . .	31
B.5	fabiasp . . . . .	37
B.6	mfsc . . . . .	41
B.7	nmfdiv . . . . .	47
B.8	nmfeu . . . . .	49
B.9	nmfsc . . . . .	51
<b>C</b>	<b>Analyzing the Results of Biclustering / Matrix Factorization</b>	<b>53</b>
C.1	plot . . . . .	53
C.2	show . . . . .	57
C.3	showSelected . . . . .	58
C.4	summary . . . . .	59
C.5	extract_bic . . . . .	60
C.6	extract_plot . . . . .	63
C.7	plotBiccluster . . . . .	65
<b>D</b>	<b>Data Generation Methods</b>	<b>67</b>
D.1	make_fabi_data . . . . .	67
D.2	make_fabi_data_pos . . . . .	69
D.3	make_fabi_data_blocks . . . . .	71

---

D.4	make_fabi_data_blocks_pos . . . . .	72
<b>E</b>	<b>Utilities</b>	<b>74</b>
E.1	fabiaVersion . . . . .	74
E.2	myImagePlot . . . . .	75
E.3	nprojfunc . . . . .	75
E.4	projfunc . . . . .	76
E.5	estimateMode . . . . .	77
E.6	eqscplotLoc . . . . .	78
<b>F</b>	<b>Data Sets</b>	<b>79</b>
F.1	Breast_A . . . . .	79
F.2	DLBCL_B . . . . .	80
F.3	Multi_A . . . . .	80

## 1 Introduction

The *fabia* package is part of the Bioconductor (<http://www.bioconductor.org>) project. The package allows to extract biclusters from data sets based on a generative model according to the FABIA method (Hochreiter *et al.*, 2010). It has been designed especially for microarray data sets, but can be used for other kinds of data sets as well.

## 2 Getting Started: FABIA

First load the *fabia* package:

```
R> library(fabia) ##load the fabia package
```

The *fabia* package imports the package methods and the function `plot` from the package *graphics*.

### 2.1 Quick start : without outputs and plots

Assume your data is in the file `datafile.csv` in a matrix like format then you can try out the following steps to extract biclusters.

Examples with outputs and plot can be found at <http://www.bioinf.jku.at/software/fabia/fabia.html>).

1. Create a working directory, e.g. `c:/fabia/data` in Windows or `/home/myself/fabia/data` in Unix. Move the data file `datafile.csv` to that directory, e.g. under Unix

```
cp datafile.csv /home/myself/fabia/data/
```

or drag the file `datafile.csv` into that directory under Windows.

2. Start R and change to the working directory. Under Windows

```
R> setwd("c:/fabia/data")
```

and under Unix

```
R> setwd("/home/myself/fabia/data")
```

You can also start R in that directory under Unix.

3. Load the library:

```
R> library(fabia) ##load the fabia package
```

4. Read the data file “`datafile.csv`”:

```
R> X <- read.table("datafile.csv",header = TRUE, sep = ",")
```

5. Select the model based on the data: 5 biclusters; sparseness 0.1; 400 cycles

```
R> res <- fabia(X,400,0.1,5)
```

6. Give summary:

```
R> summary(res)
```

7. Plot some statistics:

```
R> show(res)
```

8. Plot the factorization results:

```
R> rr <- extract_plot(res,ti="FABIA")
```

9. Plot the result as a biplot:

```
R> plot(res)
```

10. Extract biclusters:

```
R> rb <- extract_bic(res)
```

11. Show information content of the biclusters:

```
R> res@avini
```

12. List bicluster 1:

```
R> rb$bic[1,]
```

13. List bicluster 2:

```
R> rb$bic[2,]
```

14. Show bicluster 3:

```
R> rb$bic[3,]
```

15. List bicluster 4:

```
R> rb$bic[4,]
```

16. List bicluster 5:

```
R> rb$bic[5,]
```

17. Plot bicluster 1:

```
R> plotBiccluster(X,unlist(rb$bic[1,5]),unlist(rb$bic[1,3]))
```

18. Plot bicluster 2:

```
R> plotBiccluster(X,unlist(rb$bic[2,5]),unlist(rb$bic[2,3]))
```

19. Plot bicluster 3:

```
R> plotBicluster(X,unlist(rb$bic[3,5]),unlist(rb$bic[3,3]))
```

20. Plot bicluster 4:

```
R> plotBicluster(X,unlist(rb$bic[4,5]),unlist(rb$bic[4,3]))
```

21. Plot bicluster 5:

```
R> plotBicluster(X,unlist(rb$bic[5,5]),unlist(rb$bic[5,3]))
```

22. List opposite bicluster 1:

```
R> rb$bicopp[1,]
```

23. Plot opposite bicluster 1:

```
R> plotBicluster(X,unlist(rb$bicopp[1,5]),unlist(rb$bicopp[1,3]))
```

## 2.2 Test on Toy Data Set

In the following, we describe how you can test the package `fabia` on a toy data set that is generated on-line.

1. generate bicluster data, where biclusters are in block format in order to obtain a better visualization of the results. 1000 observations, 100 samples, 10 biclusters:

```
R> dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,
  f2 = 5,of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,
  mean_z = 2.0,sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)
```

2. store the generated data in variables:

```
R> X <- dat[[1]]
R> Y <- dat[[2]]
```

3. perform `fabia` (sparseness by Laplace prior) to extract biclusters; 400 cycles, sparseness 0.1 (Laplace), 13 biclusters:

```
R> resToy1 <- fabia(X,400,0.1,1.0,1.0,13)
```

4. Plot some results:

```
R> rToy1 <- extract_plot(resToy1,ti="FABIA",Y=Y)
```

5. perform `fabias` (sparseness by projection) to extract biclusters; 200 cycles, sparseness 0.5 (projection), 13 biclusters:

```
R> resToy2 <- fabias(X,200,0.5,1.0,13)
```

6. Plot some results:

```
R> rToy2 <- extract_plot(resToy2,ti="FABIAS",Y=Y)
```

7. perform `fabiap` (Laplace prior then projection) to extract biclusters; 200 cycles, sparseness 0.1 (Laplace), 13 biclusters, 0.7 sparseness loading (projection), 0.7 sparseness factors (projection):

```
R> resToy3 <- fabiap(X,200,0.1,1.0,1.0,13,0.7,0.7)
```

8. Plot some results:

```
R> rToy3 <- extract_plot(resToy3,ti="FABIAP",Y=Y)
```

9. perform `mfsc` (sparse matrix factorization), 13 biclusters, 0.7 sparseness loading (projection), 0.7 sparseness factors (projection):

```
R> resToy4 <- mfsc(X,13,sL=0.7,sZ=0.7)
```

10. Plot some results:

```
R> rToy4 <- extract_plot(resToy4,ti="MFSC",Y=Y)
```

## 2.3 Demos

The package `fabia` has some demos which can be demonstrated by `fabiaDemo`.

1. demo1: toy data.

```
R> fabiaDemo()
```

Choose “1” and you get above toy data demonstration.

2. demo2: Microarray data set of (van’t Veer *et al.*, 2002) on breast cancer.

```
R> fabiaDemo()
```

Choose “2” to extract subclasses in the data set of van’t Veer as biclusters.

3. demo3: Microarray data set of (Su *et al.*, 2002) on different mammalian.

```
R> fabiaDemo()
```

Choose “3” to check whether the different mouse and human tissue types can be extracted.

4. demo4: Microarray data set of (Rosenwald *et al.*, 2002) diffuse large-B-cell lymphoma. (Hoshida *et al.*, 2007) divided the data set into three classes

- OxPhos: oxidative phosphorylation
- BCR: B-cell response
- HR: host response

```
R> fabiaDemo()
```

Choose “4” to check whether the different classes can be extracted.

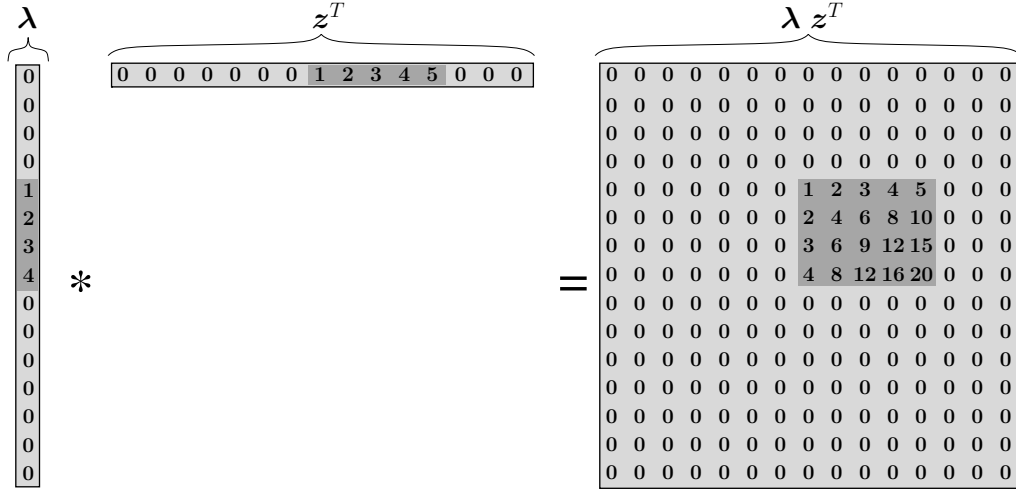


Figure 1: The outer product  $\lambda z^T$  of two sparse vectors results in a matrix with a bicluster. Note, that the non-zero entries in the vectors are adjacent to each other for visualization purposes only.

### 3 The FABIA Model

For a detailed model description see the FABIA model (Hochreiter *et al.*, 2010).

#### 3.1 Model Assumptions

We define a *bicluster* as a pair of a row (gene) set and a column (sample) set for which the rows are similar to each other on the columns and vice versa. In a multiplicative model, two vectors are similar if one is a multiple of the other, that is the angle between them is zero or as realization of random variables their correlation coefficient is one. It is clear that such a linear dependency on subsets of rows and columns can be represented as an outer product  $\lambda z^T$  of two vectors  $\lambda$  and  $z$ . The vector  $\lambda$  corresponds to a *prototype column vector* that contains zeros for genes not participating in the bicluster, whereas  $z$  is a vector of *factors* with which the prototype column vector is scaled for each sample; clearly  $z$  contains zeros for samples not participating in the bicluster. Vectors containing many zeros or values close to zero are called *sparse vectors*. Fig. 1 visualizes this representation by sparse vectors schematically.

The overall model for  $p$  biclusters and additive noise is

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon = \Lambda Z + \Upsilon, \quad (1)$$

where  $\Upsilon \in \mathbb{R}^{n \times l}$  is additive noise and  $\lambda_i \in \mathbb{R}^n$  and  $z_i \in \mathbb{R}^l$  are the sparse prototype vector and the sparse vector of factors of the  $i$ -th bicluster, respectively. The second formulation above holds if  $\Lambda \in \mathbb{R}^{n \times p}$  is the sparse prototype matrix containing the prototype vectors  $\lambda_i$  as columns and  $Z \in \mathbb{R}^{p \times l}$  is the sparse factor matrix containing the transposed factors  $z_i^T$  as rows. Note that Eq. (1) formulates biclustering as sparse matrix factorization.



According to Eq. (1), the  $j$ -th sample  $\mathbf{x}_j$ , i.e., the  $j$ -th column of  $\mathbf{X}$ , is

$$\mathbf{x}_j = \sum_{i=1}^p \lambda_i z_{ij} + \epsilon_j = \mathbf{\Lambda} \tilde{\mathbf{z}}_j + \epsilon_j, \quad (2)$$

where  $\epsilon_j$  is the  $j$ -th column of the noise matrix  $\mathbf{\Upsilon}$  and  $\tilde{\mathbf{z}}_j = (z_{1j}, \dots, z_{pj})^T$  denotes the  $j$ -th column of the matrix  $\mathbf{Z}$ . Recall that  $\mathbf{z}_i^T = (z_{i1}, \dots, z_{il})$  is the vector of values that constitutes the  $i$ -th bicluster (one value per sample), while  $\tilde{\mathbf{z}}_j$  is the vector of values that contribute to the  $j$ -th sample (one value per bicluster).

The formulation in Eq. (2) facilitates a generative interpretation by a factor analysis model with  $p$  factors

$$\mathbf{x} = \sum_{i=1}^p \lambda_i \tilde{\mathbf{z}}_i + \epsilon = \mathbf{\Lambda} \tilde{\mathbf{z}} + \epsilon, \quad (3)$$

where  $\mathbf{x}$  are the observations,  $\mathbf{\Lambda}$  is the loading matrix,  $\tilde{\mathbf{z}}_i$  is the value of the  $i$ -th factor,  $\tilde{\mathbf{z}} = (\tilde{z}_1, \dots, \tilde{z}_p)^T$  is the vector of factors, and  $\epsilon \in \mathbb{R}^n$  is the additive noise. Standard factor analysis assumes that the noise is independent of  $\tilde{\mathbf{z}}$ , that  $\tilde{\mathbf{z}}$  is  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ -distributed, and that  $\epsilon$  is  $\mathcal{N}(\mathbf{0}, \mathbf{\Psi})$ -distributed, where the covariance matrix  $\mathbf{\Psi} \in \mathbb{R}^{n \times n}$  is a diagonal matrix expressing independent Gaussian noise. The parameter  $\mathbf{\Lambda}$  explains the dependent (common) and  $\mathbf{\Psi}$  the independent variance in the observations  $\mathbf{x}$ . Normality of the additive noise in gene expression is justified by the findings in (Hochreiter *et al.*, 2006).

The unity matrix as covariance matrix for  $\tilde{\mathbf{z}}$  may be violated by overlapping biclusters, however we want to avoid to divide a real bicluster into two factors. Thus, we prefer uncorrelated factors over more sparseness. The factors can be decorrelated by setting  $\hat{\mathbf{z}} := \mathbf{A}^{-1} \tilde{\mathbf{z}}$  and  $\hat{\mathbf{\Lambda}} := \mathbf{\Lambda} \mathbf{A}$  with the symmetric invertible matrix  $\mathbf{A}^2 = \mathbb{E}(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T)$ :

$$\begin{aligned} \mathbf{\Lambda} \mathbf{z} &= \mathbf{\Lambda} \mathbf{A} \mathbf{A}^{-1} \mathbf{z} = \hat{\mathbf{\Lambda}} \hat{\mathbf{z}} \quad \text{and} \\ \mathbb{E}(\hat{\mathbf{z}} \hat{\mathbf{z}}^T) &= \mathbf{A}^{-1} \mathbb{E}(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T) \mathbf{A}^{-1} = \mathbf{A}^{-1} \mathbf{A}^2 \mathbf{A}^{-1} = \mathbf{I}. \end{aligned}$$

Standard factor analysis does not consider sparse factors and sparse loadings which are essential in our formulation to represent biclusters. Sparseness is obtained by a component-wise independent *Laplace* distribution (Hyvärinen and Oja, 1999), which is now used as a prior on the factors  $\tilde{\mathbf{z}}$  instead of the Gaussian:

$$p(\tilde{\mathbf{z}}) = \left(\frac{1}{\sqrt{2}}\right)^p \prod_{i=1}^p e^{-\sqrt{2} |\tilde{z}_i|}$$

Sparse loadings  $\lambda_i$  and, therefore sparse  $\mathbf{\Lambda}$ , are achieved by two alternative strategies. In the first model, called **FABIA**, we assume a component-wise independent *Laplace* prior for the loadings (like for the factors):

$$p(\lambda_i) = \left(\frac{1}{\sqrt{2}}\right)^n \prod_{k=1}^n e^{-\sqrt{2} |\lambda_{ki}|} \quad (4)$$

The **FABIA** model contains the product of Laplacian variables which is distributed proportionally to the 0-th order modified Bessel function of the second kind (Bithas *et al.*, 2007). For large

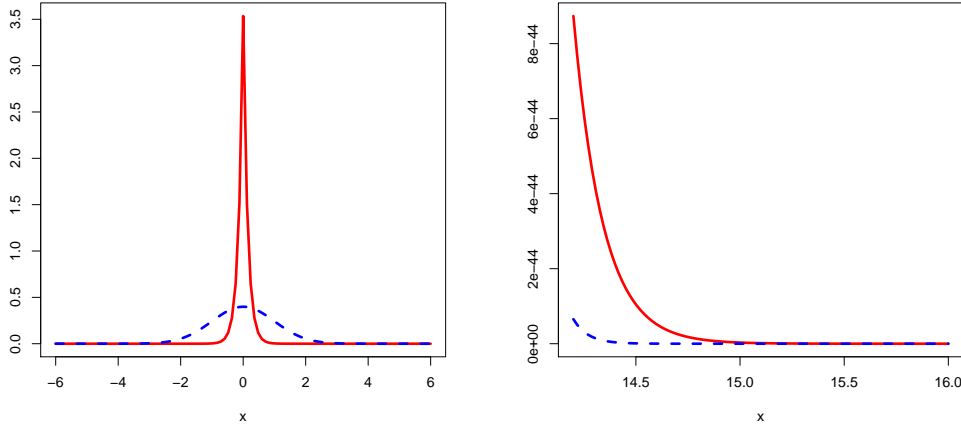


Figure 2: Left: The mode of a Laplace (red, solid) vs. a Gaussian (dashed, blue) distribution. Right: The tails of a Laplace (red, solid) vs. a Gaussian (dashed, blue) distribution

values, this Bessel function is a negative exponential function of the square root of the random variable. Therefore, the tails of the distribution are heavier than those of the Laplace distribution. The Gaussian noise, however, reduces the heaviness of the tails such that the heaviness is between Gaussian and Bessel function tails — about as heavy as the tails of the Laplacian distribution. These *heavy tails* are exactly the desired model characteristics.

The second model, called **FABIAS**, applies a prior with parameter  $\text{spL}$  on the loadings that has only support at sparse regions. Following (Hoyer, 2004), we define sparseness as

$$\text{sp}(\lambda_i) = \frac{\sqrt{n} - \sum_{k=1}^n |\lambda_{ki}| / \sum_{k=1}^n \lambda_{ki}^2}{\sqrt{n} - 1}$$

leading to the prior

$$p(\lambda_i) = \begin{cases} c & \text{for } \text{sp}(\lambda_i) \leq \text{spL} \\ 0 & \text{for } \text{sp}(\lambda_i) > \text{spL} \end{cases}. \quad (5)$$

### 3.2 Sparse Coding and Laplace Prior

Laplace prior enforces sparse codes on the factors. *Sparse coding* is the representation of items by the *strong activation* of a relatively *small set* of hidden factors while the factors are *almost constant* if not activated.

Laplace prior is suited for modeling strong activation for few samples while being otherwise almost constant. Fig. 2 left shows the Laplacian mode compared to a Gaussian mode. The Laplacian mode is higher and narrower. Fig. 2 right shows the tails of Gaussian and Laplace distribution, where the latter has higher values. This means for the Laplace distribution large values are more likely than for a Gaussian.

### 3.3 Model Selection

The free parameters  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$  can be estimated by Expectation-Maximization (EM; Dempster *et al.*, 1977). With a prior probability on the loadings, the a posteriori of the parameters is maximized like in (Hochreiter *et al.*, 2006; Talloen *et al.*, 2007).

#### 3.3.1 Variational Approach for Sparse Factors

Model selection is not straightforward because the likelihood

$$p(\mathbf{x} \mid \mathbf{\Lambda}, \mathbf{\Psi}) = \int p(\mathbf{x} \mid \tilde{\mathbf{z}}, \mathbf{\Lambda}, \mathbf{\Psi}) p(\tilde{\mathbf{z}}) d\tilde{\mathbf{z}}$$

cannot be computed analytically for a Laplacian prior  $p(\tilde{\mathbf{z}})$ . We employ a variational approach according to Girolami (2001) and Palmer *et al.* (2006) for model selection. They introduce a model family that is parametrized by  $\xi$ , where the maximum over models in this family is the true likelihood:

$$\arg \max_{\xi} p(\mathbf{x} \mid \xi) = \log p(\mathbf{x}) .$$

Using an EM algorithm, not only the likelihood with respect to the parameters  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$  is maximized, but also with respect to  $\xi$ .

In the following,  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$  denote the actual parameter estimates. According to Girolami (2001) and Palmer *et al.* (2006), we obtain

$$\begin{aligned} E(\tilde{\mathbf{z}}_j \mid \mathbf{x}_j) &= (\mathbf{\Lambda}^T \mathbf{\Psi}^{-1} \mathbf{\Lambda} + \mathbf{\Xi}_j^{-1})^{-1} \mathbf{\Lambda}^T \mathbf{\Psi}^{-1} \mathbf{x}_j \quad \text{and} \\ E(\tilde{\mathbf{z}}_j \tilde{\mathbf{z}}_j^T \mid \mathbf{x}_j) &= (\mathbf{\Lambda}^T \mathbf{\Psi}^{-1} \mathbf{\Lambda} + \mathbf{\Xi}_j^{-1})^{-1} + \\ &\quad E(\tilde{\mathbf{z}}_j \mid \mathbf{x}_j) E(\tilde{\mathbf{z}}_j \mid \mathbf{x}_j)^T , \end{aligned}$$

where  $\mathbf{\Xi}_j$  stands for  $\text{diag}(\xi_j)$ . The update for  $\xi_j$  is

$$\xi_j = \text{diag} \left( \sqrt{E(\tilde{\mathbf{z}}_j \tilde{\mathbf{z}}_j^T \mid \mathbf{x}_j)} \right) .$$

#### 3.3.2 New Update Rules for Sparse Loadings

The update rules for FABIA (Laplace prior on loadings) are

$$\begin{aligned} \mathbf{\Lambda}^{\text{new}} &= \frac{\frac{1}{l} \sum_{j=1}^l \mathbf{x}_j E(\tilde{\mathbf{z}}_j \mid \mathbf{x}_j)^T - \frac{\alpha}{l} \mathbf{\Psi} \text{sign}(\mathbf{\Lambda})}{\frac{1}{l} \sum_{j=1}^l E(\tilde{\mathbf{z}}_j \tilde{\mathbf{z}}_j^T \mid \mathbf{x}_j)} \\ \text{diag}(\mathbf{\Psi}^{\text{new}}) &= \mathbf{\Psi}^{\text{EM}} + \text{diag} \left( \frac{\alpha}{l} \mathbf{\Psi} \text{sign}(\mathbf{\Lambda}) (\mathbf{\Lambda}^{\text{new}})^T \right) \end{aligned} \tag{6}$$

where

$$\mathbf{\Psi}^{\text{EM}} = \text{diag} \left( \frac{1}{l} \sum_{j=1}^l \mathbf{x}_j \mathbf{x}_j^T - \mathbf{\Lambda}^{\text{new}} \frac{1}{l} \sum_{j=1}^l E(\tilde{\mathbf{z}}_j \mid \mathbf{x}_j) \mathbf{x}_j^T \right) .$$

The update rules for FABIAS must take into account that each  $\lambda_i$  from  $\Lambda$  has a prior with restricted support. Therefore the sparseness constraints  $\text{sp}(\lambda_i) \leq \text{spL}$  from Eq. (5) hold. These constraints are ensured by a projection of  $\lambda_i$  after each  $\Lambda$  update according to Hoyer (2004). The projection is a convex quadratic problem which minimizes the Euclidean distance to the original vector subject to the constraints. The projection problem can be solved fast by an iterative procedure where the  $l_2$ -norm of the vectors is fixed to 1. We update  $\text{diag}(\Psi^{\text{new}}) = \Psi^{\text{EM}}$  and project each updated prototype vector to a sparse vector with sparseness  $\text{spL}$  giving the overall projection:

$$\Lambda^{\text{new}} = \text{proj} \left( \frac{\frac{1}{l} \sum_{j=1}^l \mathbf{x}_j \mathbb{E}(\tilde{\mathbf{z}}_j | \mathbf{x}_j)^T}{\frac{1}{l} \sum_{j=1}^l \mathbb{E}(\tilde{\mathbf{z}}_j \tilde{\mathbf{z}}_j^T | \mathbf{x}_j)}, \text{spL} \right)$$

### 3.3.3 Extremely Sparse Priors

Some gene expression data sets are sparser than Laplacian. For example, during estimating DNA copy numbers with Affymetrix SNP 6 arrays, we observed a kurtosis larger than 30 (FABIA results shown at <http://www.bioinf.jku.at/software/fabia/fabia.html>). We want to adapt our model class to deal with such sparse data sets. Toward this end, we define extremely sparse priors both on the factors and the loadings utilizing the following (pseudo) distributions:

$$\begin{aligned} \text{Generalized Gaussians: } p(z) &\propto \exp(-|z|^\beta) \\ \text{Jeffrey's prior: } p(z) &\propto \exp(-\ln|z|) = 1/|z| \\ \text{Improper prior: } p(z) &\propto \exp(|z|^{-\beta}) \end{aligned}$$

For the first distribution, we assume  $0 < \beta \leq 1$  and for the third  $0 < \beta$ . Note, the third distribution may only exist on the interval  $[\epsilon, a]$  with  $0 < \epsilon < a$ . We assume that  $\epsilon$  is sufficiently small.

For the *loadings*, we need the derivatives of the negative log-distributions for optimizing the log-posterior. These derivatives are proportional to  $|z|^{-\text{spl}}$ , where  $\text{spl} = 0$  corresponds to the Laplace prior and  $\text{spl} > 0$  to sparser priors. The update rule is as in Eq. (6), where  $\text{sign}(\Lambda)$  is replaced by  $|\Lambda|^{-\text{spl}} \text{sign}(\Lambda)$  with element-wise operations (absolute value, sign, exponentiation, multiplication).

For the *factors*, we represent the priors through a convex variational form according to Palmer *et al.* (2006). That is possible because  $g(z) = -\ln p(\sqrt{z})$  is increasing and concave for  $z > 0$  (first order derivatives are larger and second order smaller than zero). According to Palmer *et al.* (2006), the update for  $\xi_j$  is

$$\xi_j \propto \text{diag} \left( \mathbb{E}(\tilde{\mathbf{z}}_j \tilde{\mathbf{z}}_j^T | \mathbf{x}_j)^{\text{spz}} \right)$$

for all  $\text{spz} \geq 1/2$ , where  $\text{spz} = 1/2$  ( $\beta = 1$ ) represents the Laplace prior and  $\text{spz} > 1/2$  leads to sparser priors.

## 3.4 Data Preprocessing and Initialization

The data  $\mathbf{x}$  may be centered either to zero mean or to zero median which we prefer to obtain sparser raw data. Then the data should be scaled to unit second moments to allow initialization of the parameters in the same range. See the supplementary for justification of these preprocessing steps.

The iterative model selection procedure requires initialization of the parameters  $\Lambda$ ,  $\Psi$ , and  $\xi_j$ . The simplest strategy is to initialize  $\Lambda$  randomly while ensuring that

$$\Psi = \text{diag}(\text{covar}(\mathbf{x}) - \Lambda \Lambda^T) \geq \delta > 0.$$

The variational parameter vectors  $\xi_j$  are initialized as vectors of ones. An alternative initialization strategy can be based on ICA. The ICA solution supplies factors  $\mathbf{Z}_{\text{ICA}}$  that are sparse and decorrelated.

### 3.5 Information Content of Biclusters

A highly desired property for biclustering algorithms is the ability to rank the extracted biclusters analogously to principal components which are ranked according to the data variance they explain. We rank biclusters according to the information they contain about the data. The information content of  $\tilde{z}_j$  for the  $j$ -th observation  $\mathbf{x}_j$  is the mutual information between  $\tilde{z}_j$  and  $\mathbf{x}_j$ :

$$I(\mathbf{x}_j; \tilde{z}_j) = H(\tilde{z}_j) - H(\tilde{z}_j | \mathbf{x}_j) = \frac{1}{2} \ln |\mathbf{I}_p + \Xi_j \Lambda^T \Psi^{-1} \Lambda|$$

The independence of  $\mathbf{x}_j$  and  $\tilde{z}_j$  across  $j$  gives

$$I(\mathbf{X}; \mathbf{Z}) = \frac{1}{2} \sum_{j=1}^l \ln |\mathbf{I}_p + \Xi_j \Lambda^T \Psi^{-1} \Lambda|.$$

For the FARMS summarization algorithm ( $p = 1$  and  $\Xi_j = 1$ ), this information is the negative logarithm of the I/NI call (Talloe *et al.*, 2007).

To assess the information content of one factor, we consider the case that factor  $\tilde{z}_i$  is removed from the final model. This corresponds to setting  $\xi_{ij} = 0$  (by  $\xi_{ij}$ , we denote the  $i$ -th entry in  $\xi_j$ ) and therefore the explained covariance  $\xi_{ji} \lambda_i \lambda_i^T$  is removed:

$$\mathbf{x}_j | (\tilde{z}_j \setminus z_{ij}) \sim \mathcal{N}(\Lambda \tilde{z}_j |_{z_{ij}=0}, \Psi + \xi_{ij} \lambda_i \lambda_i^T)$$

The information of  $z_{ij}$  given the other factors is

$$\begin{aligned} I(\mathbf{x}_j; z_{ij} | (\tilde{z}_j \setminus z_{ij})) &= H(z_{ij} | (\tilde{z}_j \setminus z_{ij})) - H(z_{ij} | (\tilde{z}_j \setminus z_{ij}), \mathbf{x}_j) \\ &= \frac{1}{2} \ln (1 + \xi_{ij} \lambda_i^T \Psi^{-1} \lambda_i). \end{aligned}$$

Again independence across  $j$  gives

$$I(\mathbf{X}; \mathbf{z}_i^T | (\mathbf{Z} \setminus \mathbf{z}_i^T)) = \frac{1}{2} \sum_{j=1}^l \ln (1 + \xi_{ij} \lambda_i^T \Psi^{-1} \lambda_i).$$

This information content gives that part of information in  $\mathbf{x}$  that  $\mathbf{z}_i^T$  conveys across all examples. Note that also the number of nonzero  $\lambda_i$ 's (size of the bicluster) enters the information content.

### 3.6 Extracting Members of Biclusters

After model selection in Section 3.3 and ranking of bicluster in Section 3.5, the  $i$ -th bicluster has soft gene memberships given by the absolute values of  $\lambda_i$  and soft sample memberships given by the absolute values of  $z_i^T$ .

However, applications may need hard memberships. We determine the members of bicluster  $i$  by selecting absolute values  $\lambda_{ki}$  and  $z_{ij}$  above thresholds `thresL` and `thresZ`, respectively.

First, the second moment of each factor is normalized to 1 resulting in a factor matrix  $\hat{Z}$  (in accordance with  $E(\tilde{z}\tilde{z}^T) = I$ ). Consequently,  $\Lambda$  is rescaled to  $\hat{\Lambda}$  such that  $\Lambda Z = \hat{\Lambda} \hat{Z}$ . Now the threshold `thresZ` can be chosen to determine which percentage of samples will on average belong to a bicluster. For a Laplace prior, this percentage can be computed by  $\frac{1}{2} \exp(-\sqrt{2}/\text{thresZ})$ .

In the default setting, for each factor  $\hat{z}_i$ , only one bicluster is extracted. In gene expression, an expression pattern is either absent or present but not negatively present. Therefore, the  $i$ -th bicluster is either determined by the positive or negative values of  $\hat{z}_{ij}$ . Which one of these two possibilities is chosen is decided by whether the sum over  $|\hat{z}_{ij}| > \text{thresZ}$  is larger for the positive or negative  $\hat{z}_{ij}$ .

The threshold `thresL` for the loadings is more difficult to determine, because normalization would lead to a rescaling of the already normalized factors. Since biclusters may overlap, the contribution of  $\lambda_{ki} z_{ij}$  that are relevant must be estimated. Therefore, we first estimate the standard deviation of  $\Lambda Z$  by

$$\text{sdLZ} = \sqrt{\frac{1}{p \, l \, n} \sum_{(i,j,k)=(1,1,1)}^{(p,l,n)} (\hat{\lambda}_{ki} \, \hat{z}_{ij})^2}.$$

We set this standard deviation to the product of both thresholds which is solved for `thresL`: `thresL` = `sdLZ` / `thresZ`. However, an optimal `thresL` depends on the sparseness parameters and on the characteristics of the biclustering problem.

### 3.7 C implementation of FABIA

The functions `fabia` and `fabias` are implemented in C. It turned out that these implementations are not only faster, but also more precise. Especially we use an efficient Cholesky decomposition to compute the inverse of positive definite matrices. Some R functions for computing the inverse like `solve` were inferior to that implementation.

The interface between R and C is realized by the package `Rcpp` (Samperi, 2006).

## A Class Factorization

`Factorization` is the class structure for results of matrix factorization. Especially it is designed for factor analysis used for biclustering.

The `summary` method shows information about biclusters. The `show` method plots information about biclusters. The `showSelected` method plots selected information about biclusters. The `plot` method produces biplots of biclusters.

Objects can be created by `fabia`, `fabias`, `fabiap`, `fabiasp`, `mfsc`, `nmfsc`, `nmfdiv`, and `nmfeu`.

Objects of class `Factorization` have the following slots:

1. `parameters`: Saves parameters of the factorization method in a list:
  - “method”,
  - “number of cycles”,
  - “sparseness weight”,
  - “sparseness prior for loadings”,
  - “sparseness prior for factors”,
  - “number biclusters”,
  - “projection sparseness loadings”,
  - “projection sparseness factors”,
  - “initialization range”,
  - “are loadings rescaled after each iterations”,
  - “normalization = scaling of rows”,
  - “centering method of rows”,
  - “parameter for method”.
2. `n`: number of rows, left dimension.
3. `p1`: right dimension of left matrix.
4. `p2`: left dimension of right matrix.
5. `l`: number of columns, right dimension.
6. `center`: vector of the centers.
7. `scaleData`: vector of the scaling factors.
8. `X`: centered and scaled data matrix  $n \times l$ .
9. `L`: left matrix  $n \times p1$ .
10. `Z`: right matrix  $p2 \times l$ .
11. `M`: middle matrix  $p1 \times p2$ .
12. `LZ`: matrix  $\Lambda M Z$ .
13. `U`: noise matrix.
14. `avini`: information of each bicluster, vector of length `p2`.
15. `xavini`: information extracted from each sample, vector of length `l`.
16. `ini`: information of each bicluster in each sample, matrix  $p2 \times l$ .
17. `Psi`: noise variance per row, vector of length `n`.
18. `lapla`: prior information for each sample, vector of length `l`.

This class contains the result of different matrix factorization methods. The methods may be generative or not.

Methods may be “singular value decomposition” ( $\mathbf{M}$  contains singular values as well as  $\mathbf{avini}$ ,  $\mathbf{\Lambda}$  and  $\mathbf{Z}$  are orthonormal matrices), “independent component analysis” ( $\mathbf{Z}$  contains the projection/sources,  $\mathbf{\Lambda}$  is the mixing matrix,  $\mathbf{M}$  is unity), “factor analysis” ( $\mathbf{Z}$  contains factors,  $\mathbf{\Lambda}$  the loadings,  $\mathbf{M}$  is unity,  $\mathbf{U}$  the noise,  $\mathbf{\Psi}$  the noise covariance,  $\text{lapla}$  is a variational parameter for non-Gaussian factors,  $\text{avini}$  and  $\text{ini}$  are the information the factors convey about the observations).

## B Biclustering and Matrix Factorization Methods

### B.1 fabi

Factor Analysis for Bicluster Acquisition: Laplace Prior (FABI) (Hochreiter *et al.*, 2010).

R implementation of `fabia`, therefore it is *slow*.

1. Usage: `fabi(X,p=5,alpha=0.1,cyc=500,spl=1.0,spz=1.0,center=2,norm=1)`

2. Arguments:

- `X`: the data matrix.
- `p`: number of hidden factors = number of biclusters; default = 5.
- `alpha`: sparseness loadings (0.1 - 1.0); default = 0.1.
- `cyc`: number of iterations; default = 500.
- `spl`: sparseness prior loadings (0.5 - 4.0); default = 1.0.
- `spz`: sparseness factors (0.5 - 4.0); default = 1.0.
- `center`: data centering: 1 (mean), 2 (median), > 2 (mode), 0 (no); default = 2.
- `norm`: data normalization: 1 (0.75-0.25 quantile), >1 (var=1), 0 (no); default = 1.

3. Return Values:

- object of the class `Factorization`. Containing `LZ` (estimated noise free data  $\mathbf{\Lambda Z}$ ), `L` (loadings  $\mathbf{\Lambda}$ ), `Z` (factors  $\mathbf{Z}$ ), `U` (noise  $\mathbf{X} - \mathbf{\Lambda Z}$ ), `center` (centering vector), `scaleData` (scaling vector), `X` (centered and scaled data  $\mathbf{\tilde{X}}$ ), `Psi` (noise variance  $\mathbf{\Psi}$ ), `lapla` (variational parameter), `avini` (the information which the factor  $z_{ij}$  contains about  $x_j$  averaged over  $j$ ) `xavini` (the information which the factor  $\tilde{z}_j$  contains about  $x_j$  averaged over  $j$ ) `ini` (for each  $j$  the information which the factor  $z_{ij}$  contains about  $x_j$ ).

Biclusters are found by sparse factor analysis where *both* the factors and the loadings are sparse.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T + \mathbf{\Upsilon},$$



where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{\Upsilon} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $\mathbf{x}$  that is

$$\mathbf{x} = \sum_{i=1}^p \lambda_i z_i + \epsilon = \mathbf{\Lambda} \tilde{\mathbf{z}} + \epsilon$$

The model assumptions are:

*Factor Prior is Independent Laplace:*

$$p(\tilde{\mathbf{z}}) = \left( \frac{1}{\sqrt{2}} \right)^p \prod_{i=1}^p e^{-\sqrt{2} |z_i|}$$

*Loading Prior is Independent Laplace:*

$$p(\lambda_i) = \left( \frac{1}{\sqrt{2}} \right)^n \prod_{k=1}^n e^{-\sqrt{2} |\lambda_{ki}|}$$

*Noise: Gaussian independent*

$$p(\epsilon) = \left( \frac{1}{\sqrt{2} \pi} \right)^n \prod_{k=1}^n \frac{1}{\sigma_k} e^{-\sum_{k=1}^n \frac{\epsilon_k^2}{\sigma_k^2}}$$

*Data Mean:*

$$\mathbb{E}(\mathbf{x}) = \mathbb{E}(\mathbf{\Lambda} \tilde{\mathbf{z}} + \epsilon) = \mathbf{\Lambda} \mathbb{E}(\tilde{\mathbf{z}}) + \mathbb{E}(\epsilon) = \mathbf{0}$$

*Data Covariance:*

$$\begin{aligned} \mathbb{E}(\mathbf{x} \mathbf{x}^T) &= \mathbf{\Lambda} \mathbb{E}(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T) \mathbf{\Lambda}^T + \mathbf{\Lambda} \mathbb{E}(\tilde{\mathbf{z}}) \mathbb{E}(\epsilon^T) + \mathbb{E}(\tilde{\mathbf{z}}) \mathbb{E}(\epsilon) \mathbf{\Lambda}^T + \mathbb{E}(\epsilon \epsilon^T) \\ &= \mathbf{\Lambda} \mathbf{\Lambda}^T + \text{diag}(\sigma_k^2) \end{aligned}$$

Normalizing the data to variance one for each component gives

$$\sigma_k^2 + \left( \lambda^k \right)^T \lambda^k = 1$$

Here  $\lambda^k$  is the  $k$ -th row of  $\Lambda$  (which is a row vector of length  $p$ ). We recommend to *normalize the components to variance one* in order to make the signal and noise comparable across components.

*Estimated Parameters:*  $\Lambda$  and  $\Psi$  ( $\sigma_k$ )

*Estimated Latent Variables:*  $Z$

*Estimated Noise Free Data:*  $\Lambda Z$

*Estimated Biclusters:*  $\lambda_i z_i^T$  Large values give the bicluster (ideal the nonzero values).

The model selection is performed by a variational approach according to (Girolami, 2001) and (Palmer *et al.*, 2006).

We included a prior on the parameters and minimize a lower bound on the posterior of the parameters given the data. The update of the loadings includes an additive term which pushes the loadings toward zero (Gaussian prior leads to an multiplicative factor).

The code is implemented in R , therefore it is *slow*.

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabi(X,3,0.1,20,1.0,1.0)

#-----
# DEMO1
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
```

```
resToy <- fabi(X,13,0.1,200,1.0,1.0)

rToy <- extract_plot(resToy,ti="FABI",Y=Y)

#-----
# DEMO2
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast <- fabi(X,5,0.1,200,1.0,1.0)

rBreast <- extract_plot(resBreast,ti="FABI Breast cancer(Veer)")

#sorting of predefined labels
CBreast%*%rBreast$pmZ

#-----
# DEMO3
#-----

data(Multi_A)

X <- as.matrix(XMulti)

resMulti <- fabi(X,5,0.1,200,1.0,1.0)

rMulti <- extract_plot(resMulti,ti="FABI Multiple tissues(Su)")

#sorting of predefined labels
CMulti%*%rMulti$pmZ

#-----
# DEMO4
#-----

data(DLBCL_B)

X <- as.matrix(XDLBCL)
```

```

resDLBCL <- fabi(X,5,0.1,200,1.0,1.0)

rDLBCL <- extract_plot(resDLBCL,ti="FABI Lymphoma(Rosenwald)")

#sorting of predefined labels
CDLBCL%%rDLBCL$pmZ

```

## B.2 fabia

Factor Analysis for Bicluster Acquisition: Laplace Prior (FABIA) (Hochreiter *et al.*, 2010).

C implementation of fabia.

1. Usage: `fabia(X,p=5,alpha=0.1,cyc=500,spl=1.0,spz=1.0,random=1.0,center=2,norm=1,scale=0`

2. Arguments:

- X: the data matrix.
- p: number of hidden factors = number of biclusters; default = 5.
- alpha: sparseness loadings (0.1 - 1.0); default = 0.1.
- cyc: number of iterations; default = 500.
- spl: sparseness prior loadings (0.5 - 4.0); default = 1.0.
- spz: sparseness factors (0.5 - 4.0); default = 1.0.
- random <=0: by SVD, >0: random initialization of loadings in [-random,random]; default = 1.0.
- center: data centering: 1 (mean), 2 (median), > 2 (mode), 0 (no); default = 2.
- norm: data normalization: 1 (0.75-0.25 quantile), >1 (var=1), 0 (no); default = 1.
- scale: loading vectors are scaled in each iteration to the given variance. 0.0 indicates non scaling; default = 0.0.
- lap: minimal value of the variational parameter, default = 1.

3. Return values:

- object of the class `Factorization`. Containing LZ (estimated noise free data  $\Lambda Z$ ), L (loadings  $\Lambda$ ), Z (factors  $Z$ ), U (noise  $X - \Lambda Z$ ), center (centering vector), scaleData (scaling vector), X (centered and scaled data  $X$ ), Psi (noise variance  $\Psi$ ), lapla (variational parameter), avini (the information which the factor  $z_{ij}$  contains about  $x_j$  averaged over  $j$ ) xavini (the information which the factor  $\tilde{z}_j$  contains about  $x_j$  averaged over  $j$ ) ini (for each  $j$  the information which the factor  $z_{ij}$  contains about  $x_j$ ).

Biclusters are found by sparse factor analysis where *both* the factors and the loadings are sparse.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T + \mathbf{\Upsilon},$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{\Upsilon} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $\mathbf{z}_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $\mathbf{x}$  that is

$$\mathbf{x} = \sum_{i=1}^p \lambda_i z_i + \epsilon = \mathbf{\Lambda} \tilde{\mathbf{z}} + \epsilon$$

The model assumptions are:

*Factor Prior is Independent Laplace:*

$$p(\tilde{\mathbf{z}}) = \left( \frac{1}{\sqrt{2}} \right)^p \prod_{i=1}^p e^{-\sqrt{2} |z_i|}$$

*Loading Prior is Independent Laplace:*

$$p(\lambda_i) = \left( \frac{1}{\sqrt{2}} \right)^n \prod_{k=1}^n e^{-\sqrt{2} |\lambda_{ki}|}$$

*Noise: Gaussian independent*

$$p(\epsilon) = \left( \frac{1}{\sqrt{2} \pi} \right)^n \prod_{k=1}^n \frac{1}{\sigma_k} e^{-\sum_{k=1}^n \frac{\epsilon_k^2}{\sigma_k^2}}$$

*Data Mean:*

$$\mathbb{E}(\mathbf{x}) = \mathbb{E}(\mathbf{\Lambda} \tilde{\mathbf{z}} + \epsilon) = \mathbf{\Lambda} \mathbb{E}(\tilde{\mathbf{z}}) + \mathbb{E}(\epsilon) = \mathbf{0}$$

*Data Covariance:*

$$\begin{aligned} \mathbb{E}(\mathbf{x} \mathbf{x}^T) &= \mathbf{\Lambda} \mathbb{E}(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T) \mathbf{\Lambda}^T + \mathbf{\Lambda} \mathbb{E}(\tilde{\mathbf{z}}) \mathbb{E}(\epsilon^T) + \mathbb{E}(\tilde{\mathbf{z}}) \mathbb{E}(\epsilon) \mathbf{\Lambda}^T + \mathbb{E}(\epsilon \epsilon^T) \\ &= \mathbf{\Lambda} \mathbf{\Lambda}^T + \text{diag}(\sigma_k^2) \end{aligned}$$

Normalizing the data to variance one for each component gives

$$\sigma_k^2 + \left(\boldsymbol{\lambda}^k\right)^T \boldsymbol{\lambda}^k = 1$$

Here  $\boldsymbol{\lambda}^k$  is the  $k$ -th row of  $\mathbf{\Lambda}$  (which is a row vector of length  $p$ ). We recommend to *normalize the components to variance one* in order to make the signal and noise comparable across components.

*Estimated Parameters:*  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$  ( $\sigma_k$ )

*Estimated Latent Variables:*  $\mathbf{Z}$

*Estimated Noise Free Data:*  $\mathbf{\Lambda} \mathbf{Z}$

*Estimated Biclusters:*  $\boldsymbol{\lambda}_i \mathbf{z}_i^T$  Large values give the bicluster (ideal the nonzero values).

The model selection is performed by a variational approach according to (Girolami, 2001) and (Palmer *et al.*, 2006).

We included a prior on the parameters and minimize a lower bound on the posterior of the parameters given the data. The update of the loadings includes an additive term which pushes the loadings toward zero (Gaussian prior leads to an multiplicative factor).

The code is implemented in C using the Rcpp package.

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabia(X,3,0.1,50,1.0,1.0)

#-----
# DEMO1: Toy Data
#-----

n = 1000
l= 100
p = 10
```

```

dat <- make_fabi_data_blocks(n = n,l= 1,p = p,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
ZC <- dat[[3]]
LC <- dat[[4]]

gclab <- rep.int(0,l)
gllab <- rep.int(0,n)
clab <- as.character(1:l)
llab <- as.character(1:n)
for (i in 1:p){
  for (j in ZC[i]){
    clab[j] <- paste(as.character(i),"_",clab[j],sep="")
  }
  for (j in LC[i]){
    llab[j] <- paste(as.character(i),"_",llab[j],sep="")
  }
  gclab[unlist(ZC[i])] <- gclab[unlist(ZC[i])] + p^i
  gllab[unlist(LC[i])] <- gllab[unlist(LC[i])] + p^i
}

groups <- gclab

#### FABIA

resToy1 <- fabia(X,13,0.1,400,1.0,1.0)

rToy1 <- extract_plot(resToy1,ti="FABIA",Y=Y)

raToy1 <- extract_bic(resToy1)

plotBicluster(raToy1,1)
plotBicluster(raToy1,2)
plotBicluster(raToy1,3)
plotBicluster(raToy1,4)

colnames(resToy1@X) <- clab

```

```

rownames(resToy1@X) <- l1ab

plot(resToy1,dim=c(1,2),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy1,dim=c(1,3),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy1,dim=c(2,3),label.tol=0.1,col.group = groups,lab.size=0.6)

#-----
# DEMO2: Laura van' Veer's gene expression
#       data set for breast cancer
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast1 <- fabia(X,5,0.1,400,1.0,1.0)

rBreast1 <- extract_plot(resBreast1,ti="FABIA Breast cancer(Veer)")

raBreast1 <- extract_bic(resBreast1)

plotBicluster(raBreast1,1)
plotBicluster(raBreast1,2)
plotBicluster(raBreast1,3)
plotBicluster(raBreast1,4)

plot(resBreast1,dim=c(1,2),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(1,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(1,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(1,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(2,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(2,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(2,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(3,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(3,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast1,dim=c(4,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)

#-----
# DEMO3: Su's multiple tissue types
#       gene expression data set

```



```
#-----

data(Multi_A)

X <- as.matrix(XMulti)

resMulti1 <- fabia(X,5,0.1,300,1.0,1.0)

rMulti1 <- extract_plot(resMulti1,ti="FABIA Multiple tissues(Su)")

raMulti1 <- extract_bic(resMulti1)

plotBicluster(raMulti1,1)
plotBicluster(raMulti1,2)
plotBicluster(raMulti1,3)
plotBicluster(raMulti1,4)

plot(resMulti1,dim=c(1,2),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(1,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(1,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(1,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(2,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(2,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(2,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(3,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(3,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti1,dim=c(4,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)

#-----
# DEMO4: Rosenwald's diffuse large-B-cell
#       lymphoma gene expression data set
#-----

data(DLBCL_B)

X <- as.matrix(XDLBCL)

resDLBCL1 <- fabia(X,5,0.1,400,1.0,1.0)

rDLBCL1 <- extract_plot(resDLBCL1,ti="FABIA Lymphoma(Rosenwald)")

raDLBCL1 <- extract_bic(resDLBCL1)
```

```

plotBicluster(raDLBCL1,1)
plotBicluster(raDLBCL1,2)
plotBicluster(raDLBCL1,3)
plotBicluster(raDLBCL1,4)

plot(resDLBCL1,dim=c(1,2),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(1,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(1,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(1,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(2,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(2,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(2,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(3,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(3,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL1,dim=c(4,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)

```

### B.3 fabiap

Factor Analysis for Bicluster Acquisition: Post-Projection (FABIAP) (Hochreiter *et al.*, 2010).

1. Usage: `fabia(X,p=5,alpha=0.1,cyc=500,spl=1.0,spz=1.0,sL=0.6,sZ=0.6,random=1.0,center=2,norm=1,scale=0.0,lap=1.0)`
2. Arguments:
  - `X`: the data matrix.
  - `p`: number of hidden factors = number of biclusters; default = 5.
  - `alpha`: sparseness loadings (0.1 - 1.0); default = 0.1.
  - `cyc`: number of iterations; default = 500.
  - `spl`: sparseness prior loadings (0.5 - 4.0); default = 1.0.
  - `spz`: sparseness factors (0.5 - 4.0); default = 1.0.
  - `sL`: final sparseness loadings; default = 0.6.
  - `sZ`: final sparseness factors; default = 0.6.
  - `random` ≤ 0: by SVD, > 0: random initialization of loadings in [-random,random]; default = 1.0.
  - `center`: data centering: 1 (mean), 2 (median), > 2 (mode), 0 (no); default = 2.
  - `norm`: data normalization: 1 (0.75-0.25 quantile), > 1 (var=1), 0 (no); default = 1.
  - `scale`: loading vectors are scaled in each iteration to the given variance. 0.0 indicates non scaling; default = 0.0.
  - `lap`: minimal value of the variational parameter, default = 1.
3. Return Values:

- object of the class `Factorization`. Containing `LZ` (estimated noise free data  $\Lambda Z$ ), `L` (loadings  $\Lambda$ ), `Z` (factors  $Z$ ), `U` (noise  $X - \Lambda Z$ ), `center` (centering vector), `scaledData` (scaling vector), `X` (centered and scaled data  $\tilde{X}$ ), `Psi` (noise variance  $\Psi$ ), `lapla` (variational parameter), `avini` (the information which the factor  $z_{ij}$  contains about  $x_j$  averaged over  $j$ ) `xavini` (the information which the factor  $\tilde{z}_j$  contains about  $x_j$  averaged over  $j$ ) `ini` (for each  $j$  the information which the factor  $z_{ij}$  contains about  $x_j$ ).

Biclusters are found by sparse factor analysis where *both* the factors and the loadings are sparse. Post-processing by projecting the final results to a given sparseness criterion.

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$X = \Lambda Z + \Upsilon.$$

Here  $X \in \mathbb{R}^{n \times l}$  and  $\Upsilon \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\Lambda \in \mathbb{R}^{n \times p}$ ,  $Z \in \mathbb{R}^{p \times l}$ .

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $x$  that is

$$x = \sum_{i=1}^p \lambda_i z_i + \epsilon = \Lambda \tilde{z} + \epsilon$$

The model assumptions are:

*Factor Prior is Independent Laplace:*

$$p(\tilde{z}) = \left( \frac{1}{\sqrt{2}} \right)^p \prod_{i=1}^p e^{-\sqrt{2} |z_i|}$$

*Loading Prior is Independent Laplace:*

$$p(\lambda_i) = \left( \frac{1}{\sqrt{2}} \right)^n \prod_{k=1}^n e^{-\sqrt{2} |\lambda_{ki}|}$$

*Noise: Gaussian independent*

$$p(\epsilon) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \prod_{k=1}^n \frac{1}{\sigma_k} e^{-\sum_{k=1}^n \frac{\epsilon_k^2}{\sigma_k^2}}$$

*Data Mean:*

$$E(\mathbf{x}) = E(\mathbf{\Lambda} \tilde{\mathbf{z}} + \boldsymbol{\epsilon}) = \mathbf{\Lambda} E(\tilde{\mathbf{z}}) + E(\boldsymbol{\epsilon}) = \mathbf{0}$$

*Data Covariance:*

$$\begin{aligned} E(\mathbf{x} \mathbf{x}^T) &= \mathbf{\Lambda} E(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T) \mathbf{\Lambda}^T + \mathbf{\Lambda} E(\tilde{\mathbf{z}}) E(\boldsymbol{\epsilon}^T) + E(\tilde{\mathbf{z}}) E(\boldsymbol{\epsilon}) \mathbf{\Lambda}^T + E(\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T) \\ &= \mathbf{\Lambda} \mathbf{\Lambda}^T + \text{diag}(\sigma_k^2) \end{aligned}$$

Normalizing the data to variance one for each component gives

$$\sigma_k^2 + \left( \boldsymbol{\lambda}^k \right)^T \boldsymbol{\lambda}^k = 1$$

Here  $\boldsymbol{\lambda}^k$  is the  $k$ -th row of  $\mathbf{\Lambda}$  (which is a row vector of length  $p$ ). We recommend to *normalize the components to variance one* in order to make the signal and noise comparable across components.

*Estimated Parameters:*  $\mathbf{\Lambda}$  and  $\boldsymbol{\Psi}$  ( $\sigma_k$ )

*Estimated Latent Variables:*  $\mathbf{Z}$

*Estimated Noise Free Data:*  $\mathbf{\Lambda} \mathbf{Z}$

*Estimated Biclusters:*  $\boldsymbol{\lambda}_i \mathbf{z}_i^T$  Large values give the bicluster (ideal the nonzero values).

The model selection is performed by a variational approach according to (Girolami, 2001) and (Palmer *et al.*, 2006).

We included a prior on the parameters and minimize a lower bound on the posterior of the parameters given the data. The update of the loadings includes an additive term which pushes the loadings toward zero (Gaussian prior leads to an multiplicative factor).

*Post-processing:* The final results of the loadings and the factors are projected to a sparse vector according to Hoyer, 2004: given an  $l_1$ -norm and an  $l_2$ -norm minimize the Euclidean distance to the original vector (currently the  $l_2$ -norm is fixed to 1). The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero. Instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

The code is implemented in C using the Rcpp package. The projection is implemented in R .

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100, l = 50, p = 3, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
```

```

sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabiap(X,3,0.1,50,1.0,1.0,0.6,0.6)

#-----
# DEMO1: Toy Data
#-----

n = 1000
l= 100
p = 10

dat <- make_fabi_data_blocks(n = n,l= l,p = p,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
ZC <- dat[[3]]
LC <- dat[[4]]

gclab <- rep.int(0,l)
gllab <- rep.int(0,n)
clab <- as.character(1:l)
llab <- as.character(1:n)
for (i in 1:p){
  for (j in ZC[i]){
    clab[j] <- paste(as.character(i),"_",clab[j],sep="")
  }
  for (j in LC[i]){
    llab[j] <- paste(as.character(i),"_",llab[j],sep="")
  }
  gclab[unlist(ZC[i])] <- gclab[unlist(ZC[i])] + p^i
  gllab[unlist(LC[i])] <- gllab[unlist(LC[i])] + p^i
}

groups <- gclab

#### FABIAP

```

```

resToy3 <- fabiap(X,13,0.1,400,1.0,1.0,0.6,0.6)

rToy3 <- extract_plot(resToy3,ti="FABIAP",Y=Y)

raToy3 <- extract_bic(resToy3)

plotBicluster(raToy3,1)
plotBicluster(raToy3,2)
plotBicluster(raToy3,3)
plotBicluster(raToy3,4)

colnames(resToy3@X) <- clab

rownames(resToy3@X) <- llab

plot(resToy3,dim=c(1,2),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy3,dim=c(1,3),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy3,dim=c(2,3),label.tol=0.1,col.group = groups,lab.size=0.6)

#-----
# DEMO2: Laura van' Veer's gene expression
#       data set for breast cancer
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast3 <- fabiap(X,5,0.1,400,1.0,1.0,0.6,0.6)

rBreast3 <- extract_plot(resBreast3,ti="FABIAP Breast cancer(Veer)")

raBreast3 <- extract_bic(resBreast3)

plotBicluster(raBreast3,1)
plotBicluster(raBreast3,2)
plotBicluster(raBreast3,3)
plotBicluster(raBreast3,4)

plot(resBreast3,dim=c(1,2),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(1,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(1,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(1,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)

```

```

plot(resBreast3,dim=c(2,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(2,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(2,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(3,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(3,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast3,dim=c(4,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)

#-----
# DEM03: Su's multiple tissue types
#       gene expression data set
#-----

data(Multi_A)

X <- as.matrix(XMulti)

resMulti3 <- fabiap(X,5,0.1,300,1.0,1.0,0.6,0.6)

rMulti3 <- extract_plot(resMulti3,ti="FABIAP Multiple tissues(Su)")

raMulti3 <- extract_bic(resMulti3)

plotBicluster(raMulti3,1)
plotBicluster(raMulti3,2)
plotBicluster(raMulti3,3)
plotBicluster(raMulti3,4)

plot(resMulti3,dim=c(1,2),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(1,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(1,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(1,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(2,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(2,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(2,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(3,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(3,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti3,dim=c(4,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)

#-----
# DEM04: Rosenwald's diffuse large-B-cell
#       lymphoma gene expression data set
#-----

data(DLBCL_B)

```

```

X <- as.matrix(XDLBCL)

resDLBCL3 <- fabiap(X,5,0.1,400,1.0,1.0,0.6,0.6)

rDLBCL3 <- extract_plot(resDLBCL3,ti="FABIAP Lymphoma(Rosenwald)")
raDLBCL3 <- extract_bic(resDLBCL3)

plotBicluster(raDLBCL3,1)
plotBicluster(raDLBCL3,2)
plotBicluster(raDLBCL3,3)
plotBicluster(raDLBCL3,4)

plot(resDLBCL3,dim=c(1,2),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(1,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(1,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(1,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(2,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(2,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(2,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(3,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(3,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL3,dim=c(4,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)

```

## B.4 fabias

Factor Analysis for Bicluster Acquisition: Sparseness Projection (FABIAS) (Hochreiter *et al.*, 2010).

C implementation of fabias.

1. Usage: `fabias(X,p=5,alpha=0.6,cyc=500,spz=1.0,random=1.0,center=2,norm=1,lap=1.0)`
2. Arguments:
  - X: the data matrix.
  - p: number of hidden factors = number of biclusters; default = 5.
  - alpha: sparseness loadings (0.1 - 1.0); default = 0.1.
  - cyc: number of iterations; default = 500.
  - spz: sparseness factors (0.5 - 4.0); default = 1.0.
  - random ≤ 0: by SVD, > 0: random initialization of loadings in [-random,random]; default = 1.0.
  - center: data centering: 1 (mean), 2 (median), > 2 (mode), 0 (no); default = 2.



- norm: data normalization: 1 (0.75-0.25 quantile), >1 (var=1), 0 (no); default = 1.
- lap: minimal value of the variational parameter, default = 1.

### 3. Return Values:

- object of the class `Factorization`. Containing LZ (estimated noise free data  $\Lambda Z$ ), L (loadings  $\Lambda$ ), Z (factors  $Z$ ), U (noise  $X - \Lambda Z$ ), center (centering vector), scaleData (scaling vector), X (centered and scaled data  $\tilde{X}$ ), Psi (noise variance  $\Psi$ ), lapla (variational parameter), avini (the information which the factor  $z_{ij}$  contains about  $x_j$  averaged over  $j$ ) xavini (the information which the factor  $\tilde{z}_j$  contains about  $x_j$  averaged over  $j$ ) ini (for each  $j$  the information which the factor  $z_{ij}$  contains about  $x_j$ ).

Biclusters are found by sparse factor analysis where *both* the factors and the loadings are sparse.

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$X = \Lambda Z + \Upsilon.$$

Here  $X \in \mathbb{R}^{n \times l}$  and  $\Upsilon \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\Lambda \in \mathbb{R}^{n \times p}$ ,  $Z \in \mathbb{R}^{p \times l}$ .

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $x$  that is

$$x = \sum_{i=1}^p \lambda_i z_i + \epsilon = \Lambda \tilde{z} + \epsilon$$

The model assumptions are:

*Factor Prior is Independent Laplace:*

$$p(\tilde{z}) = \left( \frac{1}{\sqrt{2}} \right)^p \prod_{i=1}^p e^{-\sqrt{2} |z_i|}$$

*Loading Prior has Finite Support:*

$$p(\lambda_i) = c \quad \text{for } \|\lambda_i\|_1 \leq k$$

$$p(\lambda_i) = 0 \quad \text{for } \|\lambda_i\|_1 > k$$

*Noise: Gaussian independent*

$$p(\epsilon) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \prod_{k=1}^n \frac{1}{\sigma_k} e^{-\sum_{k=1}^n \frac{\epsilon_k^2}{\sigma_k^2}}$$

*Data Mean:*

$$\mathbf{E}(\mathbf{x}) = \mathbf{E}(\mathbf{\Lambda} \tilde{\mathbf{z}} + \epsilon) = \mathbf{\Lambda} \mathbf{E}(\tilde{\mathbf{z}}) + \mathbf{E}(\epsilon) = \mathbf{0}$$

*Data Covariance:*

$$\begin{aligned} \mathbf{E}(\mathbf{x} \mathbf{x}^T) &= \mathbf{\Lambda} \mathbf{E}(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T) \mathbf{\Lambda}^T + \mathbf{\Lambda} \mathbf{E}(\tilde{\mathbf{z}}) \mathbf{E}(\epsilon^T) + \mathbf{E}(\tilde{\mathbf{z}}) \mathbf{E}(\epsilon) \mathbf{\Lambda}^T + \mathbf{E}(\epsilon \epsilon^T) \\ &= \mathbf{\Lambda} \mathbf{\Lambda}^T + \text{diag}(\sigma_k^2) \end{aligned}$$

Normalizing the data to variance one for each component gives

$$\sigma_k^2 + \left( \boldsymbol{\lambda}^k \right)^T \boldsymbol{\lambda}^k = 1$$

Here  $\boldsymbol{\lambda}^k$  is the  $k$ -th row of  $\mathbf{\Lambda}$  (which is a row vector of length  $p$ ). We recommend to *normalize the components to variance one* in order to make the signal and noise comparable across components.

*Estimated Parameters:*  $\mathbf{\Lambda}$  and  $\sigma_k$

*Estimated Latent Variables:*  $\mathbf{Z}$

*Estimated Noise Free Data:*  $\mathbf{\Lambda} \mathbf{Z}$

*Estimated Biclusters:*  $\boldsymbol{\lambda}_i \mathbf{z}_i^T$  Larges values give the bicluster (ideal the nonzero values).

The model selection is performed by a variational approach according to (Girolami, 2001) and (Palmer *et al.*, 2006).

The prior has finite support, therefore after each update of the loadings they are projected to the finite support. The projection is done according to (Hoyer, 2004): given an  $l_1$ -norm and an  $l_2$ -norm minimize the Euclidean distance to the original vector (currently the  $l_2$ -norm is fixed to 1). The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero. Instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

$$\text{sparseness}(\boldsymbol{\lambda}_i) = \frac{\sqrt{n} - \sum_{k=1}^n |\lambda_{ki}|}{\sqrt{n} - 1} / \sum_{k=1}^n \lambda_{ki}^2$$

The code is implemented in C using the Rcpp package.

EXAMPLE:

```

#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabias(X,3,0.6,50,1.0)

#-----
# DEMO1: Toy Data
#-----

n = 1000
l= 100
p = 10

dat <- make_fabi_data_blocks(n = n,l= l,p = p,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
ZC <- dat[[3]]
LC <- dat[[4]]

gclab <- rep.int(0,l)
gllab <- rep.int(0,n)
clab <- as.character(1:l)
llab <- as.character(1:n)
for (i in 1:p){
  for (j in ZC[i]){
    clab[j] <- paste(as.character(i),"_",clab[j],sep="")
  }
  for (j in LC[i]){
    llab[j] <- paste(as.character(i),"_",llab[j],sep="")
  }
  gclab[unlist(ZC[i])] <- gclab[unlist(ZC[i])] + p^i
}

```

```

  gllab[unlist(LC[i])] <- gllab[unlist(LC[i])] + p^i
}

groups <- gclab

#### FABIAS

resToy2 <- fabias(X,13,0.6,400,1.0)

rToy2 <- extract_plot(resToy2,ti="FABIAS",Y=Y)

raToy2 <- extract_bic(resToy2)

plotBicluster(raToy2,1)
plotBicluster(raToy2,2)
plotBicluster(raToy2,3)
plotBicluster(raToy2,4)

colnames(resToy2@X) <- clab

rownames(resToy2@X) <- llab

plot(resToy2,dim=c(1,2),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy2,dim=c(1,3),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy2,dim=c(2,3),label.tol=0.1,col.group = groups,lab.size=0.6)

#-----
# DEMO2: Laura van' Veer's gene expression
#       data set for breast cancer
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast2 <- fabias(X,5,0.6,300,1.0)

rBreast2 <- extract_plot(resBreast2,ti="FABIAS Breast cancer(Veer)")

raBreast2 <- extract_bic(resBreast2)

plotBicluster(raBreast2,1)

```

```

plotBicluster(raBreast2,2)
plotBicluster(raBreast2,3)
plotBicluster(raBreast2,4)

plot(resBreast2,dim=c(1,2),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(1,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(1,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(1,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(2,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(2,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(2,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(3,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(3,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast2,dim=c(4,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)

#-----
# DEM03: Su's multiple tissue types
#       gene expression data set
#-----

data(Multi_A)

X <- as.matrix(XMulti)

resMulti2 <- fabias(X,5,0.6,300,1.0)

rMulti2 <- extract_plot(resMulti2,ti="FABIAS Multiple tissues(Su)")

raMulti2 <- extract_bic(resMulti2)

plotBicluster(raMulti2,1)
plotBicluster(raMulti2,2)
plotBicluster(raMulti2,3)
plotBicluster(raMulti2,4)

plot(resMulti2,dim=c(1,2),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(1,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(1,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(1,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(2,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(2,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(2,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(3,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(3,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti2,dim=c(4,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)

```

```

#-----
# DEMO4: Rosenwald's diffuse large-B-cell
#       lymphoma gene expression data set
#-----

data(DLBCL_B)

X <- as.matrix(XDLBCL)

resDLBCL2 <- fabias(X,5,0.6,300,1.0)

rDLBCL2 <- extract_plot(resDLBCL2,ti="FABIAS Lymphoma(Rosenwald)")

raDLBCL2 <- extract_bic(resDLBCL2)

plotBicluster(raDLBCL2,1)
plotBicluster(raDLBCL2,2)
plotBicluster(raDLBCL2,3)
plotBicluster(raDLBCL2,4)

plot(resDLBCL2,dim=c(1,2),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(1,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(1,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(1,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(2,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(2,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(2,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(3,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(3,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL2,dim=c(4,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)

```

## B.5 fabiasp

Factor Analysis for Bicluster Acquisition: Sparseness Projection (FABIASP) (Hochreiter *et al.*, 2010).

R implementation of fabias, therefore it is *slow*.

1. Usage: `fabiasp(X,p=5,alpha=0.6,cyc=500,spz=1.0,center=2,norm=1)`
2. Arguments:

- $X$ : the data matrix.
- $p$ : number of hidden factors = number of biclusters; default = 5.
- $\alpha$ : sparseness loadings (0.1 - 1.0); default = 0.1.
- $cyc$ : number of iterations; default = 500.
- $spz$ : sparseness factors (0.5 - 4.0); default = 1.0.
- $center$ : data centering: 1 (mean), 2 (median), > 2 (mode), 0 (no); default = 2.
- $norm$ : data normalization: 1 (0.75-0.25 quantile), >1 (var=1), 0 (no); default = 1.

### 3. Return Values:

- object of the class `Factorization`. Containing  $LZ$  (estimated noise free data  $\Lambda Z$ ),  $L$  (loadings  $\Lambda$ ),  $Z$  (factors  $Z$ ),  $U$  (noise  $X - \Lambda Z$ ),  $center$  (centering vector),  $scaleData$  (scaling vector),  $X$  (centered and scaled data  $\tilde{X}$ ),  $Psi$  (noise variance  $\Psi$ ),  $lapla$  (variational parameter),  $avini$  (the information which the factor  $z_{ij}$  contains about  $x_j$  averaged over  $j$ )  $xavini$  (the information which the factor  $\tilde{z}_j$  contains about  $x_j$  averaged over  $j$ )  $ini$  (for each  $j$  the information which the factor  $z_{ij}$  contains about  $x_j$ ).

Biclusters are found by sparse factor analysis where *both* the factors and the loadings are sparse.

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$X = \Lambda Z + \Upsilon.$$

Here  $X \in \mathbb{R}^{n \times l}$  and  $\Upsilon \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\Lambda \in \mathbb{R}^{n \times p}$ ,  $Z \in \mathbb{R}^{p \times l}$ .

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $x$  that is

$$x = \sum_{i=1}^p \lambda_i z_i + \epsilon = \Lambda \tilde{z} + \epsilon$$

The model assumptions are:

*Factor Prior is Independent Laplace:*

$$p(z) = \left( \frac{1}{\sqrt{2}} \right)^p \prod_{i=1}^p e^{-\sqrt{2} |z_i|}$$

*Loading Prior has Finite Support:*

$$p(\boldsymbol{\lambda}_i) = c \quad \text{for } \|\boldsymbol{\lambda}_i\|_1 \leq k$$

$$p(\boldsymbol{\lambda}_i) = 0 \quad \text{for } \|\boldsymbol{\lambda}_i\|_1 > k$$

*Noise: Gaussian independent*

$$p(\boldsymbol{\epsilon}) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \prod_{k=1}^n \frac{1}{\sigma_k} e^{-\sum_{k=1}^n \frac{\epsilon_k^2}{\sigma_k^2}}$$

*Data Mean:*

$$\mathbf{E}(\mathbf{x}) = \mathbf{E}(\mathbf{\Lambda} \tilde{\mathbf{z}} + \boldsymbol{\epsilon}) = \mathbf{\Lambda} \mathbf{E}(\tilde{\mathbf{z}}) + \mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$$

*Data Covariance:*

$$\begin{aligned} \mathbf{E}(\mathbf{x} \mathbf{x}^T) &= \mathbf{\Lambda} \mathbf{E}(\tilde{\mathbf{z}} \tilde{\mathbf{z}}^T) \mathbf{\Lambda}^T + \mathbf{\Lambda} \mathbf{E}(\tilde{\mathbf{z}}) \mathbf{E}(\boldsymbol{\epsilon}^T) + \mathbf{E}(\tilde{\mathbf{z}}) \mathbf{E}(\boldsymbol{\epsilon}) \mathbf{\Lambda}^T + \mathbf{E}(\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T) \\ &= \mathbf{\Lambda} \mathbf{\Lambda}^T + \text{diag}(\sigma_k^2) \end{aligned}$$

Normalizing the data to variance one for each component gives

$$\sigma_k^2 + \left( \boldsymbol{\lambda}^k \right)^T \boldsymbol{\lambda}^k = 1$$

Here  $\boldsymbol{\lambda}^k$  is the  $k$ -th row of  $\mathbf{\Lambda}$  (which is a row vector of length  $p$ ). We recommend to *normalize the components to variance one* in order to make the signal and noise comparable across components.

*Estimated Parameters:  $\mathbf{\Lambda}$  and  $\boldsymbol{\Psi}$  ( $\sigma_k$ )*

*Estimated Latent Variables:  $\mathbf{Z}$*

*Estimated Noise Free Data:  $\mathbf{\Lambda} \mathbf{Z}$*

*Estimated Biclusters:  $\boldsymbol{\lambda}_i \mathbf{z}_i^T$*  Larges values give the bicluster (ideal the nonzero values).

The model selection is performed by a variational approach according to (Girolami, 2001) and (Palmer *et al.*, 2006).

The prior has finite support, therefore after each update of the loadings they are projected to the finite support. The projection is done according to (Hoyer, 2004): given an  $l_1$ -norm and an  $l_2$ -norm minimize the Euclidean distance to the original vector (currently the  $l_2$ -norm is fixed to 1). The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero. Instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

$$\text{sparseness}(\boldsymbol{\lambda}_i) = \frac{\sqrt{n} - \sum_{k=1}^n |\lambda_{ki}|}{\sqrt{n} - 1} \bigg/ \sum_{k=1}^n \lambda_{ki}^2$$



The code is implemented in R , therefore it is *slow*.

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabiasp(X,3,0.6,50,1.0)

\dontrun{
#-----
# DEMO1
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resToy <- fabiasp(X,13,0.6,200,1.0)

rToy <- extract_plot(resToy,ti="FABIASP",Y=Y)

#-----
# DEMO2
#-----

data(Breast_A)
```

```
X <- as.matrix(XBreast)

resBreast <- fabiasp(X,5,0.6,200,1.0)

rBreast <- extract_plot(resBreast,ti="FABIASP Breast cancer(Veer)")

#sorting of predefined labels
CBreast%%rBreast$pmZ

#-----
# DEMO3
#-----

data(Multi_A)

X <- as.matrix(XMulti)

resMulti <- fabiasp(X,5,0.6,200,1.0)

rMulti <- extract_plot(resMulti,"ti=FABIASP Multiple tissues(Su)")

#sorting of predefined labels
CMulti%%rMulti$pmZ

#-----
# DEMO4
#-----

data(DLBCL_B)

X <- as.matrix(XDLBCL)

resDLBCL <- fabiasp(X,5,0.6,200,1.0)

rDLBCL <- extract_plot(resDLBCL,ti="FABIASP Lymphoma(Rosenwald)")

#sorting of predefined labels
CDLBCL%%rDLBCL$pmZ
```

## B.6 mfsc

Sparse Matrix Factorization for bicluster analysis (MFSC) (Hochreiter *et al.*, 2010).

1. Usage: `mfsc(X,p=5,cyc=100,sL=0.6,sZ=0.6,center=2,norm=1)`
2. Arguments:
  - `X`: the data matrix.
  - `p`: number of hidden factors = number of biclusters; default = 5.
  - `cyc`: maximal number of iterations; default = 100.
  - `sL`: final sparseness loadings; default = 0.6.
  - `sZ`: final sparseness factors; default = 0.6.
  - `center`: data centering: 1 (mean), 2 (median), > 2 (mode), 0 (no); default = 2.
  - `norm`: data normalization: 1 (0.75-0.25 quantile), >1 (var=1), 0 (no); default = 1.
3. Return Values:
  - object of the class `Factorization`. Containing `LZ` (estimated noise free data  $\Lambda Z$ ), `L` (loadings  $\Lambda$ ), `Z` (factors  $Z$ ), `U` (noise  $X - \Lambda Z$ ), `center` (centering vector), `scaleData` (scaling vector), `X` (centered and scaled data  $X$ ).

Biclusters are found by sparse matrix factorization where *both* factors are sparse.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $\mathbf{z}_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

*No noise assumption:* In contrast to factor analysis there is no noise assumption.

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $\mathbf{x}$  that is

$$\mathbf{x} = \sum_{i=1}^p \lambda_i \mathbf{z}_i = \mathbf{\Lambda} \tilde{\mathbf{z}}$$

*Estimated Parameters:*  $\mathbf{\Lambda}$  and  $\mathbf{Z}$

*Estimated Biclusters:*  $\lambda_i \mathbf{z}_i^T$  Largest values give the bicluster (ideal the nonzero values).

The model selection is performed by a constraint optimization according to (Hoyer, 2004). The Euclidean distance (the Frobenius norm) is minimized subject to sparseness constraints:

$$\min_{\mathbf{\Lambda}, \mathbf{Z}} \|\mathbf{X} - \mathbf{\Lambda} \mathbf{Z}\|_F^2$$

$$\text{subject to } \|\mathbf{\Lambda}\|_F^2 = 1$$

$$\text{subject to } \|\mathbf{\Lambda}\|_1 = k_L$$

$$\text{subject to } \|\mathbf{Z}\|_F^2 = 1$$

$$\text{subject to } \|\mathbf{Z}\|_1 = k_Z$$

Model selection is done by gradient descent on the Euclidean objective and thereafter projection of single vectors of  $\mathbf{\Lambda}$  and single vectors of  $\mathbf{Z}$  to fulfill the sparseness constraints.

The projection minimize the Euclidean distance to the original vector given an  $l_1$ -norm and an  $l_2$ -norm.

The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero. Instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

$$\text{sparseness}(\lambda_i) = \frac{\sqrt{n} - \sum_{k=1}^n |\lambda_{ki}| / \sum_{k=1}^n \lambda_{ki}^2}{\sqrt{n} - 1}$$

The code is implemented in R .

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100, l = 50, p = 3, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- mfsc(X, 3, 30, 0.6, 0.6)
```

```

#-----
# DEM01: Toy Data
#-----

n = 1000
l= 100
p = 10

dat <- make_fabi_data_blocks(n = n,l= l,p = p,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
ZC <- dat[[3]]
LC <- dat[[4]]

gclab <- rep.int(0,l)
gllab <- rep.int(0,n)
clab <- as.character(1:l)
llab <- as.character(1:n)
for (i in 1:p){
  for (j in ZC[i]){
    clab[j] <- paste(as.character(i),"_",clab[j],sep="")
  }
  for (j in LC[i]){
    llab[j] <- paste(as.character(i),"_",llab[j],sep="")
  }
  gclab[unlist(ZC[i])] <- gclab[unlist(ZC[i])] + p^i
  gllab[unlist(LC[i])] <- gllab[unlist(LC[i])] + p^i
}

groups <- gclab

#### MFSC

resToy4 <- mfsc(X,13,100,0.6,0.6)

rToy4 <- extract_plot(resToy4,ti="MFSC",Y=Y)

raToy4 <- extract_bic(resToy4,thresZ=0.01,thresL=0.05)

plotBicluster(raToy4,1)

```

```

plotBicluster(raToy4,2)
plotBicluster(raToy4,3)
plotBicluster(raToy4,4)

colnames(resToy4@X) <- clab

rownames(resToy4@X) <- llab

plot(resToy4,dim=c(1,2),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy4,dim=c(1,3),label.tol=0.1,col.group = groups,lab.size=0.6)
plot(resToy4,dim=c(2,3),label.tol=0.1,col.group = groups,lab.size=0.6)

#-----
# DEM02: Laura van' Veer's gene expression
#       data set for breast cancer
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast4 <- mfsc(X,5,100,0.6,0.6)

rBreast4 <- extract_plot(resBreast4,ti="MFSC Breast cancer(Veer)")

raBreast4 <- extract_bic(resBreast4,thresZ=0.01,thresL=0.05)

plotBicluster(raBreast4,1)
plotBicluster(raBreast4,2)
plotBicluster(raBreast4,3)
plotBicluster(raBreast4,4)

plot(resBreast4,dim=c(1,2),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(1,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(1,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(1,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(2,3),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(2,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(2,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(3,4),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(3,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)
plot(resBreast4,dim=c(4,5),label.tol=0.03,col.group=CBreast,lab.size=0.6)

```

```

#-----
# DEMO3: Su's multiple tissue types
#       gene expression data set
#-----

data(Multi_A)

X <- as.matrix(XMulti)

resMulti4 <- mfsc(X,5,100,0.6,0.6)

rMulti4 <- extract_plot(resMulti4,ti="MFSC Multiple tissues(Su)")

raMulti4 <- extract_bic(resMulti4,thresZ=0.01,thresL=0.05)

plotBicluster(raMulti4,1)
plotBicluster(raMulti4,2)
plotBicluster(raMulti4,3)
plotBicluster(raMulti4,4)

plot(resMulti4,dim=c(1,2),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(1,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(1,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(1,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(2,3),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(2,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(2,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(3,4),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(3,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)
plot(resMulti4,dim=c(4,5),label.tol=0.01,col.group=CMulti,lab.size=0.6)

#-----
# DEMO4: Rosenwald's diffuse large-B-cell
#       lymphoma gene expression data set
#-----

data(DLBCL_B)

X <- as.matrix(XDLBCL)

resDLBCL4 <- mfsc(X,5,100,0.6,0.6)

rDLBCL4 <- extract_plot(resDLBCL4,ti="MFSC Lymphoma(Rosenwald)")

```

```

raDLBCL4 <- extract_bic(resDLBCL4,thresZ=0.01,thresL=0.05)

plotBicluster(raDLBCL4,1)
plotBicluster(raDLBCL4,2)
plotBicluster(raDLBCL4,3)
plotBicluster(raDLBCL4,4)

plot(resDLBCL4,dim=c(1,2),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(1,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(1,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(1,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(2,3),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(2,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(2,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(3,4),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(3,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)
plot(resDLBCL4,dim=c(4,5),label.tol=0.03,col.group=CDLBCL,lab.size=0.6)

```

## B.7 nmfdiv

Non-negative Matrix Factorization with Kullback-Leibler divergence as objective.

1. Usage: `nmfdiv(X,p=5,cyc=100)`
2. Arguments:
  - `X`: the data matrix.
  - `p`: number of hidden factors = number of biclusters; default = 5.
  - `cyc`: maximal number of iterations; default = 100.
3. Return Values:
  - object of the class `Factorization`. Containing `LZ` (estimated noise free data  $\Lambda Z$ ), `L` (loadings  $\Lambda$ ), `Z` (factors  $Z$ ), `U` (noise  $X - \Lambda Z$ ), `X` (scaled data  $X$ ).

$$X = \Lambda Z$$

$$X = \sum_{i=1}^p \lambda_i z_i^T$$



*Estimated Parameters:  $\Lambda$  and  $Z$*

The model selection is performed according to (Lee and Seung, 1999, 2001).

*objective:*

$$D(\mathbf{A} \parallel \mathbf{B}) = \sum_{ij} \left( A_{ij} \log \frac{A_{ij}}{B_{ij}} + A_{ij} - B_{ij} \right)$$

*update:*

$$L_{ik} = L_{ik} \frac{\sum_{j=1}^n Z_{ji} V_{jk} / (\Lambda Z)_{jk}}{\sum_{j=1}^n Z_{ji}}$$

$$Z_{ji} = Z_{ji} \frac{\sum_{k=1}^l L_{ik} V_{jk} / (\Lambda Z)_{jk}}{\sum_{k=1}^l L_{ik}}$$

or in matrix notation with “\*” and “/” as element-wise operators:

$$\Lambda = \Lambda * ((\mathbf{X} / (\Lambda \mathbf{Z})) t(\mathbf{Z})) / \text{rowSums}(\mathbf{Z})$$

$$\mathbf{Z} = \mathbf{Z} * (t(\Lambda) (\mathbf{X} / (\Lambda \mathbf{Z}))) / \text{colSums}(\Lambda)$$

The code is implemented in R .

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
X <- abs(X)

resEx <- nmfddiv(X,3)
```

```

#-----
# DEMO
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
X <- abs(X)

resToy <- nmfddiv(X,13)

rToy <- extract_plot(resToy,ti="NMFDIV",Y=Y)

```

## B.8 nmfeu

Non-negative Matrix Factorization with Euclidean distance as objective.

1. Usage: `nmfeu(X,p=5,cyc=100)`
2. Arguments:
  - `X`: the data matrix.
  - `p`: number of hidden factors = number of biclusters; default = 5.
  - `cyc`: maximal number of iterations; default = 100.
3. Return Values:
  - object of the class `Factorization`. Containing `LZ` (estimated noise free data  $\Lambda Z$ ), `L` (loadings  $\Lambda$ ), `Z` (factors  $Z$ ), `U` (noise  $X - \Lambda Z$ ), `X` (scaled data  $X$ ).

$$X = \Lambda Z$$

$$X = \sum_{i=1}^p \lambda_i z_i^T$$

*Estimated Parameters:  $\Lambda$  and  $Z$*

The model selection is performed according to (Lee and Seung, 2001; Paatero and Tapper, 1997).

*objective:*

$$\|\mathbf{A} - \mathbf{B}\|_F^2 = \sum_{ij} (A_{ij} - B_{ij})^2$$

*update:*

$$L_{ik} = L_{ik} \frac{(\mathbf{\Lambda}^T \mathbf{X})_{ik}}{(\mathbf{\Lambda}^T \mathbf{\Lambda} \mathbf{Z})_{ik}}$$

$$Z_{ji} = Z_{ji} \frac{(\mathbf{X} \mathbf{Z}^T)_{ji}}{(\mathbf{\Lambda} \mathbf{Z} \mathbf{Z}^T)_{ji}}$$

or in matrix notation with “\*” and “/” as element-wise operators:

$$\mathbf{Z} = \mathbf{Z} * (t(\mathbf{\Lambda}) \mathbf{X}) / (t(\mathbf{\Lambda}) \mathbf{\Lambda} \mathbf{Z})$$

$$\mathbf{\Lambda} = \mathbf{\Lambda} * (\mathbf{X} t(\mathbf{Z})) / (\mathbf{\Lambda} \mathbf{Z} t(\mathbf{Z}))$$

The code is implemented in R .

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
X <- abs(X)

resEx <- nmfeu(X,3)

#-----
```

```
# DEMO
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
X <- abs(X)

resToy <- nmfeu(X,13)

rToy <- extract_plot(resToy,ti="NMFEU",Y=Y)
```

## B.9 nmfsc

Non-negative Sparse Matrix Factorization with sparseness constraints.

1. Usage: `nmfsc(X,p=5,cyc=100,sL=0.6,sZ=0.6)`
2. Arguments:
  - `X`: the data matrix.
  - `p`: number of hidden factors = number of biclusters; default = 5.
  - `cyc`: maximal number of iterations; default = 100.
  - `sL`: sparseness loadings; default = 0.6.
  - `sZ`: sparseness factors; default = 0.6.
3. Return Values:
  - object of the class `Factorization`. Containing `LZ` (estimated noise free data  $\Lambda Z$ ), `L` (loadings  $\Lambda$ ), `Z` (factors  $Z$ ), `U` (noise  $X - \Lambda Z$ ), `X` (data  $X$ ).

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$X = \Lambda Z.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$ ,  $\boldsymbol{\lambda}_i \in \mathbb{R}^n$ ,  $\mathbf{z}_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

If the nonzero components of the sparse vectors are grouped together then the outer product results in a matrix with a nonzero block and zeros elsewhere.

For a single data vector  $\mathbf{x}$  that is

$$\mathbf{x} = \sum_{i=1}^p \boldsymbol{\lambda}_i z_i = \mathbf{\Lambda} \mathbf{z}$$

*Estimated Parameters:*  $\mathbf{\Lambda}$  and  $\mathbf{Z}$

*Estimated Biclusters:*  $\boldsymbol{\lambda}_i \mathbf{z}_i^T$  Large values give the bicluster (ideal the nonzero values).

The model selection is performed by a constraint optimization according to (Hoyer, 2004). The Euclidean distance (the Frobenius norm) is minimized subject to sparseness and non-negativity constraints:

$$\min_{\mathbf{\Lambda}, \mathbf{Z}} \|\mathbf{x} - \mathbf{\Lambda} \mathbf{Z}\|_F^2$$

$$\text{subject to } \|\mathbf{\Lambda}\|_F^2 = 1$$

$$\text{subject to } \|\mathbf{\Lambda}\|_1 = k_L$$

$$\text{subject to } \mathbf{\Lambda} \geq \mathbf{0}$$

$$\text{subject to } \|\mathbf{Z}\|_F^2 = 1$$

$$\text{subject to } \|\mathbf{Z}\|_1 = k_Z$$

$$\text{subject to } \mathbf{Z} \geq \mathbf{0}$$

Model selection is done by gradient descent on the Euclidean objective and thereafter projection of single vectors of  $\mathbf{\Lambda}$  and single vectors of  $\mathbf{Z}$  to fulfill the sparseness and non-negativity constraints.

The projection minimize the Euclidean distance to the original vector given an  $l_1$ -norm and an  $l_2$ -norm and enforcing non-negativity.

The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero. Instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

$$\text{sparseness}(\boldsymbol{\lambda}_i) = \frac{\sqrt{n} - \sum_{k=1}^n |\lambda_{ki}| / \sum_{k=1}^n \lambda_{ki}^2}{\sqrt{n} - 1}$$

The code is implemented in R .

EXAMPLE:

```

#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
X <- abs(X)

resEx <- nmfsc(X,3,30,0.6,0.6)

#-----
# DEMO
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]
X <- abs(X)

resToy <- nmfsc(X,13,100,0.6,0.6)

rToy <- extract_plot(resToy,ti="NMFSC",Y=Y)

```

## C Analyzing the Results of Biclustering / Matrix Factorization

### C.1 plot

Plot of a Matrix Factorization.

Produces a (biplot) of a Matrix Factorization result.

1. Usage: `plot(x,y, ...)`
2. Arguments:
  - `x`: object of the class `Factorization`.

- `y`: missing, not used.
- `...`: with following parameters
  - `Rm`: row weighting vector.
  - `Cm`: column weighting vector.
  - `dim`: optional principal factors that are plotted along the horizontal and vertical axis. Defaults to `c(1,2)`.
  - `zoom`: optional zoom factor for row and column items. Defaults to `c(1,1)`.
  - `col.group`: optional vector (character or numeric) indicating the different groupings of the columns. Defaults to 1.
  - `colors`: vector specifying the colors for the annotation of the plot; the first two elements concern the rows; the third till the last element concern the columns; the first element will be used to color the unlabeled rows; the second element for the labeled rows and the remaining elements to give different colors to different groups of columns.
  - `col.areas`: logical value indicating whether columns should be plotted as squares with areas proportional to their marginal mean and colors representing the different groups (TRUE), or with symbols representing the groupings and identical size (FALSE). Defaults to TRUE.
  - `col.symbols`: vector of symbols when `col.areas=FALSE` corresponds to the `pch` argument of the function `plot`.
  - `sampleNames`: Either a logical vector of length one or a character vector of length equal to the number of samples in the dataset. If a logical is provided, sample names will be displayed on the plot (TRUE; default) or not (FALSE); if a character vector is provided, the names provided will be used to label the samples instead of the default column names.
  - `rot`: rotation of plot. Defaults to `c(-1,-1)`.
  - `labels`: character vector to be used for labeling points on the graph; if NULL, the row names of `x` are used instead.
  - `label.tol`: numerical value specifying either the percentile (`label.tol<=1`) of rows or the number of rows (`label.tol>1`) most distant from the plot-center (0,0) that are labeled and are plotted as circles with area proportional to the marginal means of the original data.
  - `lab.size`: size of identifying labels for row- and column-items as `cex` parameter of the `text` function.
  - `col.size`: size in mm of the column symbols.
  - `row.size`: size in mm of the row symbols.
  - `do.smoothScatter`: use `smoothScatter` or not instead of plotting individual points.
  - `do.plot`: produce a plot or not.
  - `...`: further arguments to `eqscaleplotLoc` which draws the canvas for the plot; useful for adding a main or a custom sub.

and with the default values

- `Rm=rep(1,nrow(X))`,
- `Cm=rep(1,ncol(X))`,

- `dim = c(1, 2),`
- `zoom = rep(1, 2),`
- `col.group = rep(1, ncol(X)),`
- `colors = c("orange1", "red", rainbow(length(unique(col.group))), start=2/6, end=4/6),`
- `col.areas = TRUE,`
- `col.symbols = c(1, rep(2, length(unique(col.group))))),`
- `sampleNames = TRUE,`
- `rot = rep(-1, length(dim)),`
- `labels = NULL,`
- `label.tol = 1,`
- `lab.size = 0.725,`
- `col.size = 10,`
- `row.size = 10,`
- `do.smoothScatter = FALSE,`
- `do.plot = TRUE,`

### 3. Return Values:

- Rows: a list with the X and Y coordinates of the rows and an indication `Select` of whether the row was selected according to `label.tol`.
- Columns: a list with the X and Y coordinates of the columns.

The function `plot.fabia` is based on the function `plot.mpm` in the R package `mpm` (Version: 1.0-16, Date: 2009-08-26, Title: Multivariate Projection Methods, Maintainer: Tobias Verbeke <tobias.verbeke@openanalytics.be>, Author: Luc Wouters <wouters\_luc@telenet.be>).

Biclusters are found by sparse factor analysis where *both* the factors and the loadings are sparse.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T + \mathbf{\Upsilon},$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{\Upsilon} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $\mathbf{z}_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

For noise free projection like independent component analysis we set the noise term to zero:  $\mathbf{\Upsilon} = \mathbf{0}$ .

The argument `label.tol` can be used to select the most informative rows, i.e. rows that are most distant from the center of the plot (smaller 1: percentage of rows, larger 1: number of rows).

Only these row-items are then labeled and represented as circles with their areas proportional to the row weighting.



If the column-items are grouped these groups can be visualized by colors given by `col.group`.  
Implementation in R .

EXAMPLE:

```
n=200
l=100
p=4

dat <- make_fabi_data_blocks(n = n,l= l,p = p,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
ZC <- dat[[3]]
LC <- dat[[4]]

resEx <- fabia(X,p,0.1,400,1.0,1.0)

gclab <- rep.int(0,l)
gllab <- rep.int(0,n)
clab <- as.character(1:l)
llab <- as.character(1:n)
for (i in 1:p){
  for (j in ZC[i]){
    clab[j] <- paste(as.character(i),"_",clab[j],sep="")
  }
  for (j in LC[i]){
    llab[j] <- paste(as.character(i),"_",llab[j],sep="")
  }
  gclab[unlist(ZC[i])] <- gclab[unlist(ZC[i])] + p^i
  gllab[unlist(LC[i])] <- gllab[unlist(LC[i])] + p^i
}

labels <- paste(as.character(clab) , "_",as.character(1:l),sep="")

groups <- gclab

colnames(resEx@X) <- clab
```

```

rlabels <- paste(as.character(llab) , "_", as.character(1:1), sep="")

rownames(resEx@X) <- llab

plot(resEx, dim=c(1,2), label.tol=0.1, col.group = groups, lab.size=0.6)

plot(resEx, dim=c(1,3), label.tol=0.1, col.group = groups, lab.size=0.6)
plot(resEx, dim=c(2,3), label.tol=0.1, col.group = groups, lab.size=0.6)

```

## C.2 show

Plots Statistics of a Matrix Factorization.

This function plots statistics on a matrix factorization which is stored as an instance of `Factorization-class`.

1. Usage: `show(object)`
2. Arguments:
  - `object`: An instance of `Factorization-class`.
3. Return Values:
  - no value.

Following is plotted:

- 1) the information content of biclusters.
- 2) the information content of samples.
- 3) the loadings per bicluster.
- 4) the factors per bicluster.

Implementation in R .

EXAMPLE:

```

#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100, l= 50, p = 3, f1 = 5, f2 = 5,

```

```
of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,  
sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)  
  
X <- dat[[1]]  
  
resEx <- fabia(X, 3, 0.1, 100, 1.0, 1.0)  
  
show(resEx)
```

### C.3 showSelected

Plots selected Statistics of a Matrix Factorization.

This function plots selected statistics on a matrix factorization which is stored as an instance of `Factorization-class`.

1. Usage: `showSelected(object, which=c(1,2,3,4))`
2. Arguments:
  - `object`: An instance of `Factorization-class`.
  - `which`: A list of which plots should be generated: 1=the information content of biclusters, 2=the information content of samples, 3=the loadings per bicluster, 4=the factors per bicluster, default `c(1,2,3,4)`.
3. Return Values:
  - no value.

Following is plotted depending on the plot selection variable "which":

- 1) the information content of biclusters.
- 2) the information content of samples.
- 3) the loadings per bicluster.
- 4) the factors per bicluster.

Implementation in R .

EXAMPLE:

```

#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]

resEx <- fabia(X,3,0.1,100,1.0,1.0)

showSelected(resEx,which=1)
showSelected(resEx,which=2)

```

## C.4 summary

Summary of Matrix Factorization.

This function gives information on a matrix factorization which is stored as an instance of `Factorization-class`.

1. Usage: `summary(object, ...)`
2. Arguments:
  - `object`: An instance of `Factorization-class`.
  - `...`: further arguments.
3. Return Values:
  - no value.

First, it gives the number of rows and columns of the original matrix.

Then the number of clusters for rows and columns is given.

Then for the row cluster the information content is given.

Then for each column its information is given.

Then for each column cluster a summary is given.

Then for each row cluster a summary is given.

Implementation in R .

EXAMPLE:

```

#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]

resEx <- fabia(X,3,0.1,100,1.0,1.0)

summary(resEx)

```

### C.5 extract\_bic

Extraction of Biclusters.

1. Usage: `extract_bic(fact,thresZ=0.5,thresL=NULL)`
2. Arguments:
  - `fact`: object of class `Factorization`.
  - `thresZ`: threshold for sample belonging to bicluster; default 0.5.
  - `thresL`: threshold for loading belonging to bicluster (if not given it is estimated).
3. Return Values:
  - `bic`: extracted biclusters.
  - `numn`: indexes for the extracted biclusters.
  - `bicopp`: extracted opposite biclusters.
  - `numnopp`: indexes for the extracted opposite biclusters.
  - `X`: scaled and centered data matrix.
  - `np`: number of biclusters.

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{\Upsilon} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

$\mathbf{\Upsilon}$  is the Gaussian noise with a diagonal covariance matrix which entries are given by  $\Psi_i$ .

The  $\mathbf{Z}$  is locally approximated by a Gaussian with inverse variance given by  $\text{lapla}$ .

The  $\lambda_i$  and  $z_i$  are used to extract the bicluster  $i$ , where a threshold determines which observations and which samples belong to the bicluster.

In `bic` the biclusters are extracted according to the largest absolute values of the component  $i$ , i.e. the largest values of  $\lambda_i$  and the largest values of  $z_i$ . The factors  $z_i$  are normalized to variance 1.

The components of `bic` are `binp`, `bixv`, `bixn`, `biypv`, and `biypn`.

`binp` give the size of the bicluster: number observations and number samples. `bixv` gives the values of the extracted observations that have absolute values above a threshold. They are sorted. `bixn` gives the extracted observation names (e.g. gene names). `biypv` gives the values of the extracted samples that have absolute values above a threshold. They are sorted. `biypn` gives the names of the extracted samples (e.g. sample names).

In `bicopp` the opposite cluster to the biclusters are given. Opposite means that the negative pattern is present.

The components of opposite clusters `bicopp` are `binp`, `bixv`, `bixn`, `biypnv`, and `biypnn`.

`binp` give the size of the opposite bicluster: number observations and number samples. `bixv` gives the values of the extracted observations that have absolute values above a threshold. They are sorted. `bixn` gives the extracted observation names (e.g. gene names). `biypnv` gives the values of the opposite extracted samples that have absolute values above a threshold. They are sorted. `biypnn` gives the names of the opposite extracted samples (e.g. sample names).

That means the samples are divided into two groups where one group shows large positive values and the other group has negative values with large absolute values. That means a observation pattern can be switched on or switched off relative to the average value.

`numn` gives the indexes of `bic` with components: `numng = bix` and `numnp = biypn`.

`numn` gives the indexes of `bicopp` with components: `numng = bix` and `numnn = biypnn`.

Implementation in R .

EXAMPLE:

```
#-----
# TEST
#-----
```

```
dat <- make_fabi_data_blocks(n = 100, l = 50, p = 3, f1 = 5, f2 = 5,
```

```

    of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
    sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabia(X, 3, 0.1, 20, 1.0, 1.0)

rEx <- extract_bic(resEx)

rEx$bic[1,]
rEx$bic[2,]
rEx$bic[3,]

#-----
# DEMO1
#-----

dat <- make_fabi_data_blocks(n = 1000, l = 100, p = 10, f1 = 5, f2 = 5,
    of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
    sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resToy <- fabia(X, 13, 0.1, 200, 1.0, 1.0)

rToy <- extract_bic(resToy)

resToy@avini

rToy$bic[1,]
rToy$bic[2,]
rToy$bic[3,]

#-----
# DEMO2
#-----

data(Breast_A)

X <- as.matrix(XBreast)

```

```

resBreast <- fabia(X,5,0.1,200,1.0,1.0)

rBreast <- extract_bic(resBreast)

resBreast@avini

rBreast$bic[1,]
rBreast$bic[2,]
rBreast$bic[3,]

```

## C.6 extract\_plot

Plotting of Biclustering Results.

1. Usage: `extract_plot(fact,thresZ=0.5,ti="",thresL=NULL,Y=NULL,which=c(1,2,3,4,5,6,7,8))`
2. Arguments:
  - `fact`: object of class `Factorization`.
  - `thresZ`: threshold for sample belonging to bicluster; default 0.5.
  - `thresL`: threshold for loading belonging to bicluster (estimated if not given).
  - `ti`: plot title; default "".
  - `Y`: noise free data matrix.
  - `which`: which plot is shown: 1=noise free data (if available), 2=data, 3=reconstructed data, 4=error, 5=absolute factors, 6=absolute loadings, 7=reconstructed matrix sorted, 8=original matrix sorted; default `c(1,2,3,4,5,6,7,8)`.
3. Return Values:
  - `biclust`: clusters of kmeans clustering.
  - `pmZ`: permutation matrix of  $z$  from kmeans clustering.
  - `pmL`: permutation matrix of  $L$  from kmeans clustering.
  - `nL`: normalized loadings (left matrix).
  - `nZ`: normalized factors (right matrix).
  - `Xord`: sorted original matrix according to kmeans on  $Z$  and kmeans on  $L$ .
4. The method produces following plots depending what plots are chosen by the "which" variable:
  - $Y$ : noise free data (if available),
  - $X$ : data,
  - $\Lambda Z$ : reconstructed data,
  - $\Lambda Z - X$ : error,



- `abs(Z)`: absolute factors,
- `abs()`: absolute loadings,
- `pmLZpmZ`: reconstructed matrix sorted,
- `pmLXpmZ`: original matrix sorted.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T + \mathbf{\Upsilon},$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{\Upsilon} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $\mathbf{z}_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

The  $\lambda_i$  and  $\mathbf{z}_i^T$  are used to extract the bicluster  $i$ , where a threshold determines which observations and which samples belong to the bicluster.

The method produces the plots given above at “Plots”.

For sorting first `kmeans` is performed on the  $p$  dimensional space and then the vectors which belong to the same cluster are put together.

This sorting is in general not able to visualize all biclusters as blocks if they overlap.

The `kmeans` clusters are given by `biclust` with components `biclustx` (the clustered observations) and `biclusty` (the clustered samples).

Implementation in R .

EXAMPLE:

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100, l = 50, p = 3, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabia(X, 3, 0.3, 20, 1.0, 1.0)
```

```

rEx <- extract_plot(resEx,ti="FABIA",Y=Y)

#-----
# DEMO1
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resToy <- fabia(X,13,0.1,200,1.0,1.0)

rToy <- extract_plot(resToy,ti="FABIA",Y=Y)

#-----
# DEMO2
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast <- fabia(X,5,0.1,200,1.0,1.0)

rBreast <- extract_plot(resBreast,ti="FABIA Breast cancer(Veer)")

#sorting of predefined labels
CBreast%%rBreast$pmZ

```

## C.7 plotBicluster

Plots a bicluster.

1. Usage: `plotBicluster(r,p,opp=FALSE,zlim=NULL,title=NULL,which=c(1, 2))`
2. Arguments:
  - `r`: the result of `extract_bic`.

- `p`: the bicluster to plot.
- `opp`: plot opposite bicluster, default = FALSE.
- `zlim`: vector containing a low and high value to use for the color scale.
- `title`: title of the plot.
- `which`: which plots are shown: 1=data matrix with bicluster on upper left, 2=plot of the bicluster; default `c(1, 2)`.

One bicluster is visualized by two plots. The variable "which" indicates which plots should be shown.

Plot1 (which=1): The data matrix is sorted such that the bicluster appear at the upper left corner.

The bicluster is marked by a rectangle.

Plot2 (which=2): Only the bicluster is presented.

Implementation in R .

```
#-----
# TEST
#-----

dat <- make_fabi_data_blocks(n = 100,l= 50,p = 3,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

resEx <- fabia(X,3,0.1,20,1.0,1.0)

rEx <- extract_bic(resEx)

plotBicluster(rEx,p=1)

\dontrun{
#-----
# DEMO1
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)
```

```

X <- dat[[1]]
Y <- dat[[2]]

resToy <- fabia(X,13,0.1,200,1.0,1.0)

rToy <- extract_bic(resToy)

plotBicluster(rToy,p=1)

#-----
# DEMO2
#-----

data(Breast_A)

X <- as.matrix(XBreast)

resBreast <- fabia(X,5,0.1,200,1.0,1.0)

rBreast <- extract_bic(resBreast)

plotBicluster(rBreast,p=1)

```

## D Data Generation Methods

### D.1 make\_fabi\_data

Generation of bicluster data.

1. Usage: `make_fabi_data(n,l,p,f1,f2,of1,of2,sd_noise,sd_z_noise,mean_z,sd_z,sd_l_noise,mean_l,sd_l)`
2. Arguments:
  - `n`: number of observations.
  - `l`: number of samples.
  - `p`: number of biclusters.
  - `f1`:  $l/f1$  max. additional samples are active in a bicluster.
  - `f2`:  $n/f2$  max. additional observations that form a pattern in a bicluster.
  - `of1`: minimal active samples in a bicluster.

- of2: minimal observations that form a pattern in a bicluster.
- sd\_noise: Gaussian zero mean noise std on data matrix.
- sd\_z\_noise: Gaussian zero mean noise std for deactivated hidden factors.
- mean\_z: Gaussian mean for activated factors.
- sd\_z: Gaussian std for activated factors.
- sd\_l\_noise: Gaussian zero mean noise std if no observation patterns are present.
- mean\_l: Gaussian mean for observation patterns.
- sd\_l: Gaussian std for observation patterns.

3. Return values:

- $X$ : the noisy data  $X$  from  $\mathbb{R}^{n \times l}$ .
- $Y$ : the noise free data  $Y$  from  $\mathbb{R}^{n \times l}$ .
- $ZC$ : list where  $i$ th element gives samples belonging to  $i$ th bicluster.
- $LC$ : list where  $i$ th element gives observations belonging to  $i$ th bicluster.

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$X = \Lambda Z + \Upsilon.$$

Here  $X \in \mathbb{R}^{n \times l}$  and  $\Upsilon \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\Lambda \in \mathbb{R}^{n \times p}$ ,  $Z \in \mathbb{R}^{p \times l}$ .

Sequentially  $\lambda_i$  are generated using `n`, `f2`, `of2`, `sd_l_noise`, `mean_l`, `sd_l`. `of2` gives the minimal observations participating in a bicluster to which between 0 and  $n/f2$  observations are added, where the number is uniformly chosen. `sd_l_noise` gives the noise of observations not participating in the bicluster. `mean_l` and `sd_l` determines the Gaussian from which the values are drawn for the observations that participate in the bicluster. The sign of the mean is randomly chosen for each component.

Sequentially  $z_i$  are generated using `l`, `f1`, `of1`, `sd_z_noise`, `mean_z`, `sd_z`. `of1` gives the minimal samples participating in a bicluster to which between 0 and  $l/f1$  samples are added, where the number is uniformly chosen. `sd_z_noise` gives the noise of samples not participating in the bicluster. `mean_z` and `sd_z` determines the Gaussian from which the values are drawn for the samples that participate in the bicluster.

$\Upsilon$  is the overall Gaussian zero mean noise generated by `sd_noise`.

Implementation in R .

EXAMPLE:

```

#-----
# DEMO
#-----

dat <- make_fabi_data(n = 1000, l = 100, p = 10, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

myImagePlot(Y)
x11()
myImagePlot(X)

```

## D.2 make\_fabi\_data\_pos

Generation of bicluster data.

1. Usage: `make_fabi_data_pos(n, l, p, f1, f2, of1, of2, sd_noise, sd_z_noise, mean_z, sd_z, sd_l_noise, mean_l, sd_l)`
2. Arguments:
  - `n`: number of observations.
  - `l`: number of samples.
  - `p`: number of biclusters.
  - `f1`:  $l/f1$  max. additional samples are active in a bicluster.
  - `f2`:  $n/f2$  max. additional observations that form a pattern in a bicluster.
  - `of1`: minimal active samples in a bicluster.
  - `of2`: menial observations that form a pattern in a bicluster.
  - `sd_noise`: Gaussian zero mean noise std on data matrix.
  - `sd_z_noise`: Gaussian zero mean noise std for deactivated hidden factors.
  - `mean_z`: Gaussian mean for activated factors.
  - `sd_z`: Gaussian std for activated factors.
  - `sd_l_noise`: Gaussian zero mean noise std if no observation patterns are present.
  - `mean_l`: Gaussian mean for observation patterns.
  - `sd_l`: Gaussian std for observation patterns.
3. Return values:
  - `X`: the noisy data  $\mathbf{X}$  from  $\mathbb{R}^{n \times l}$ .

- $Y$ : the noise free data  $Y$  from  $\mathbb{R}^{n \times l}$ .
- $ZC$ : list where  $i$ th element gives samples belonging to  $i$ th bicluster.
- $LC$ : list where  $i$ th element gives observations belonging to  $i$ th bicluster.

Essentially the model is the sum of outer products of vectors:

$$X = \sum_{i=1}^p \lambda_i z_i^T + \Upsilon,$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$X = \Lambda Z + \Upsilon.$$

Here  $X \in \mathbb{R}^{n \times l}$  and  $\Upsilon \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\Lambda \in \mathbb{R}^{n \times p}$ ,  $Z \in \mathbb{R}^{p \times l}$ .

Sequentially  $\lambda_i$  are generated using `n`, `f2`, `of2`, `sd_l_noise`, `mean_l`, `sd_l`. `of2` gives the minimal observations participating in a bicluster to which between 0 and  $n/f2$  observations are added, where the number is uniformly chosen. `sd_l_noise` gives the noise of observations not participating in the bicluster. `mean_l` and `sd_l` determines the Gaussian from which the values are drawn for the observations that participate in the bicluster. "POS": The sign of the mean is fixed.

Sequentially  $z_i$  are generated using `l`, `f1`, `of1`, `sd_z_noise`, `mean_z`, `sd_z`. `of1` gives the minimal samples participating in a bicluster to which between 0 and  $l/f1$  samples are added, where the number is uniformly chosen. `sd_z_noise` gives the noise of samples not participating in the bicluster. `mean_z` and `sd_z` determines the Gaussian from which the values are drawn for the samples that participate in the bicluster.

$\Upsilon$  is the overall Gaussian zero mean noise generated by `sd_noise`.

Implementation in R .

EXAMPLE:

```
#-----
# DEMO
#-----

dat <- make_fabi_data_pos(n = 1000, l = 100, p = 10, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

myImagePlot(Y)
x11()
myImagePlot(X)
```

### D.3 make\_fabi\_data\_blocks

Generation of bicluster data with bicluster blocks.

1. Usage: `make_fabi_data_blocks(n,l,p,f1,f2,of1,of2,sd_noise,sd_z_noise,mean_z,sd_z,sd_l_noise,mean_l,sd_l)`

2. Arguments:

- `n`: number of observations.
- `l`: number of samples.
- `p`: number of biclusters.
- `f1`:  $l/f1$  max. additional samples are active in a bicluster.
- `f2`:  $n/f2$  max. additional observations that form a pattern in a bicluster.
- `of1`: minimal active samples in a bicluster.
- `of2`: minimal observations that form a pattern in a bicluster.
- `sd_noise`: Gaussian zero mean noise std on data matrix.
- `sd_z_noise`: Gaussian zero mean noise std for deactivated hidden factors.
- `mean_z`: Gaussian mean for activated factors.
- `sd_z`: Gaussian std for activated factors.
- `sd_l_noise`: Gaussian zero mean noise std if no observation patterns are present.
- `mean_l`: Gaussian mean for observation patterns.
- `sd_l`: Gaussian std for observation patterns.

3. Return Values:

- `X`: the noisy data  $\mathbf{X}$  from  $\mathbb{R}^{n \times l}$ .
- `Y`: the noise free data  $\mathbf{Y}$  from  $\mathbb{R}^{n \times l}$ .
- `ZC`: list where  $i$ th element gives samples belonging to  $i$ th bicluster.
- `LC`: list where  $i$ th element gives observations belonging to  $i$ th bicluster.

Bicluster data is generated for visualization because the biclusters are now in block format. That means observations and samples that belong to a bicluster are consecutive. This allows visual inspection because the user can identify blocks and whether they have been found or reconstructed.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T + \mathbf{\Upsilon},$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$



Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{Y} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}^n$ ,  $z_i \in \mathbb{R}^l$ ,  $\mathbf{A} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

Sequentially  $\lambda_i$  are generated using `n`, `f2`, `of2`, `sd_l_noise`, `mean_l`, `sd_l`. `of2` gives the minimal observations participating in a bicluster to which between 0 and  $n/f2$  observations are added, where the number is uniformly chosen. `sd_l_noise` gives the noise of observations not participating in the bicluster. `mean_l` and `sd_l` determines the Gaussian from which the values are drawn for the observations that participate in the bicluster. The sign of the mean is randomly chosen for each component.

Sequentially  $z_i$  are generated using `l`, `f1`, `of1`, `sd_z_noise`, `mean_z`, `sd_z`. `of1` gives the minimal samples participating in a bicluster to which between 0 and  $l/f1$  samples are added, where the number is uniformly chosen. `sd_z_noise` gives the noise of samples not participating in the bicluster. `mean_z` and `sd_z` determines the Gaussian from which the values are drawn for the samples that participate in the bicluster.

$\mathbf{Y}$  is the overall Gaussian zero mean noise generated by `sd_noise`.

Implementation in R .

EXAMPLE:

```
#-----
# DEMO
#-----

dat <- make_fabi_data_blocks(n = 1000, l = 100, p = 10, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

Y <- dat[[1]]
X <- dat[[2]]

myImagePlot(Y)
x11()
myImagePlot(X)
```

#### D.4 make\_fabi\_data\_blocks\_pos

Generation of bicluster data with bicluster blocks.

1. Usage: `make_fabi_data_blocks_pos(n, l, p, f1, f2, of1, of2, sd_noise,`  
`sd_z_noise, mean_z, sd_z, sd_l_noise, mean_l, sd_l)`
2. Arguments:

- $n$ : number of observations.
- $l$ : number of samples.
- $p$ : number of biclusters.
- $f1$ :  $l/f1$  max. additional samples are active in a bicluster.
- $f2$ :  $n/f2$  max. additional observations that form a pattern in a bicluster.
- $of1$ : minimal active samples in a bicluster.
- $of2$ : minimal observations that form a pattern in a bicluster.
- $sd\_noise$ : Gaussian zero mean noise std on data matrix.
- $sd\_z\_noise$ : Gaussian zero mean noise std for deactivated hidden factors.
- $mean\_z$ : Gaussian mean for activated factors.
- $sd\_z$ : Gaussian std for activated factors.
- $sd\_l\_noise$ : Gaussian zero mean noise std if no observation patterns are present.
- $mean\_l$ : Gaussian mean for observation patterns.
- $sd\_l$ : Gaussian std for observation patterns.

### 3. Return Values:

- $X$ : the noisy data  $\mathbf{X}$  from  $\mathbb{R}^{n \times l}$ .
- $Y$ : the noise free data  $\mathbf{Y}$  from  $\mathbb{R}^{n \times l}$ .
- $ZC$ : list where  $i$ th element gives samples belonging to  $i$ th bicluster.
- $LC$ : list where  $i$ th element gives observations belonging to  $i$ th bicluster.

Bicluster data is generated for visualization because the biclusters are now in block format. That means observations and samples that belong to a bicluster are consecutive. This allows visual inspection because the user can identify blocks and whether they have been found or reconstructed.

Essentially the model is the sum of outer products of vectors:

$$\mathbf{X} = \sum_{i=1}^p \lambda_i \mathbf{z}_i^T + \mathbf{\Upsilon},$$

where the number of summands  $p$  is the number of biclusters. The matrix factorization is

$$\mathbf{X} = \mathbf{\Lambda} \mathbf{Z} + \mathbf{\Upsilon}.$$

Here  $\mathbf{X} \in \mathbb{R}^{n \times l}$  and  $\mathbf{\Upsilon} \in \mathbb{R}^{n \times l}$ ,  $\lambda_i \in \mathbb{R}$ ,  $\mathbf{z}_i \in \mathbb{R}^l$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{Z} \in \mathbb{R}^{p \times l}$ .

Sequentially  $\lambda_i$  are generated using  $n$ ,  $f2$ ,  $of2$ ,  $sd\_l\_noise$ ,  $mean\_l$ ,  $sd\_l$ .  $of2$  gives the minimal observations participating in a bicluster to which between 0 and  $n/f2$  observations are added, where the number is uniformly chosen.  $sd\_l\_noise$  gives the noise of observations not participating in the bicluster.  $mean\_l$  and  $sd\_l$  determines the Gaussian from which the values are drawn for the observations that participate in the bicluster. "POS": The sign of the mean is fixed.

Sequentially  $z_i$  are generated using  $l$ ,  $f1$ ,  $of1$ ,  $sd\_z\_noise$ ,  $mean\_z$ ,  $sd\_z$ .  $of1$  gives the minimal samples participating in a bicluster to which between 0 and  $l/f1$  samples are added, where the number is uniformly chosen.  $sd\_z\_noise$  gives the noise of samples not participating in the bicluster.  $mean\_z$  and  $sd\_z$  determines the Gaussian from which the values are drawn for the samples that participate in the bicluster.

$\Upsilon$  is the overall Gaussian zero mean noise generated by  $sd\_noise$ .

Implementation in R .

EXAMPLE:

```
#-----
# DEMO
#-----

dat <- make_fabi_data_blocks_pos(n = 1000, l = 100, p = 10, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 3.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

Y <- dat[[1]]
X <- dat[[2]]

myImagePlot(Y)
x11()
myImagePlot(X)
```

## E Utilities

### E.1 fabiaVersion

Display version info for package and for FABIA.

1. Usage: `fabiaVersion()`

EXAMPLE:

```
fabiaVersion()
```

## E.2 myImagePlot

Plotting of a matrix.

1. Usage: `myImagePlot(x,xLabels=NULL, yLabels=NULL, zlim=NULL, title=NULL)`
2. Arguments:

- `x`: the matrix.
- `xLabels`: vector of strings to label the rows (default `"rownames(x)"`).
- `yLabels`: vector of strings to label the columns (default `"colnames(x)"`).
- `zlim`: vector containing a low and high value to use for the color scale.
- `title`: title of the plot.

Plotting a table of numbers as an image using R .

The color scale is based on the highest and lowest values in the matrix.

Program has been obtained by <http://www.phaget4.org/R/myImagePlot.R>

EXAMPLE:

```
#-----
# DEMO
#-----

dat <- make_fabi_data_blocks(n = 1000,l= 100,p = 10,f1 = 5,f2 = 5,
  of1 = 5,of2 = 10,sd_noise = 3.0,sd_z_noise = 0.2,mean_z = 2.0,
  sd_z = 1.0,sd_l_noise = 0.2,mean_l = 3.0,sd_l = 1.0)

X <- dat[[1]]
Y <- dat[[2]]

myImagePlot(X)
```

## E.3 nprojfunc

Projection of a vector to a sparse non-negative vector with given sparseness and given  $l_2$ -norm.

1. Usage: `nprojfunc(s, k1, k2)`
2. Arguments:
  - `s`: data vector.

- k1: sparseness,  $l_1$  norm constraint.
- k2:  $l_2$  norm constraint.

### 3. Return Values:

- v: non-negative sparse projected vector.

The projection minimize the Euclidean distance to the original vector given an  $l_1$ -norm and an  $l_2$ -norm and enforcing non-negativity.

The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero.

In the applications, instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

$$\text{sparseness}(\mathbf{v}) = \frac{\sqrt{l} - \sum_{j=1}^l |v_j| / \sum_{j=1}^l v_j^2}{\sqrt{l} - 1}$$

The code is implemented in R .

### EXAMPLE:

```
#-----
# DEMO
#-----

size <- 30
sparseness <- 0.7

s <- as.vector(rnorm(size))
sp <- sqrt(1.0*size)-(sqrt(1.0*size)-1.0)*sparseness

ss <- nprojfunc(s,k1=sp,k2=1)

s
ss
```

## E.4 projfunc

Projection of a vector to a sparse vector with given sparseness and given  $l_2$ -norm.

1. Usage: `projfunc(s, k1, k2)`
2. Arguments:

- *s*: data vector.
- *k1*: sparseness,  $l_1$  norm constraint.
- *k2*:  $l_2$  norm constraint.

### 3. Return Values:

- *v*: sparse projected vector.

The projection is done according to (Hoyer, 2004): given an  $l_1$ -norm and an  $l_2$ -norm minimize the Euclidean distance to the original vector. The projection is a convex quadratic problem which is solved iteratively where at each iteration at least one component is set to zero.

In the applications, instead of the  $l_1$ -norm a sparseness measurement is used which relates the  $l_1$ -norm to the  $l_2$ -norm:

$$\text{sparseness}(\mathbf{v}) = \frac{\sqrt{l} - \sum_{j=1}^l |v_j| / \sum_{j=1}^l v_j^2}{\sqrt{l} - 1}$$

The code is implemented in R .

#### EXAMPLE:

```
#-----
# DEMO
#-----

size <- 30
sparseness <- 0.7

s <- as.vector(rnorm(size))
sp <- sqrt(1.0*size)-(sqrt(1.0*size)-1.0)*sparseness

ss <- projfunc(s,k1=sp,k2=1)

s
ss
```

## E.5 estimateMode

Estimation of the modes of the rows of a matrix.

1. Usage: `estimateMode(X,maxiter=50,tol=0.001,alpha=0.1,a1=4.0,G1=FALSE)`
2. Arguments:

- $X$ : matrix of which the modes of the rows are estimated.
- `maxiter`: maximal number of iterations; default = 50.
- `tol`: tolerance for stopping; default = 0.001.
- `alpha`: learning rate; default = 0.1.
- `a1`: parameter of the width of the given distribution; default = 4.
- `G1`: kind of distribution, TRUE:  $\frac{1}{a_1} \ln(\cosh(a_1 x))$ , FALSE:  $-\frac{1}{a_1} \exp(-\frac{a_1}{2} x^2)$ ; default = FALSE.

### 3. Return Values:

- `u`: the vector of estimated modes.
- `xu`:  $X - u\mathbf{1}^T$  the mode centered data.

The mode is estimated by contrast functions  $G1 \frac{1}{a_1} \ln(\cosh(a_1 x))$  or  $G2 -\frac{1}{a_1} \exp(-\frac{a_1}{2} x^2)$ . The estimation is performed by gradient descent initialized by the median.

Implementation in R .

### EXAMPLE:

```
#-----
# DEMO
#-----

dat <- make_fabi_data_blocks_pos(n = 100, l = 50, p = 10, f1 = 5, f2 = 5,
  of1 = 5, of2 = 10, sd_noise = 2.0, sd_z_noise = 0.2, mean_z = 2.0,
  sd_z = 1.0, sd_l_noise = 0.2, mean_l = 3.0, sd_l = 1.0)

X <- dat[[1]]

modes <- estimateMode(X)

modes$u - apply(X, 1, median)
```

## E.6 eqscplotLoc

Local Version of eqscplot from the package MASS.

Plots with Geometrically Equal Scales: Version of a scatterplot with scales chosen to be equal on both axes, that is 1cm represents the same units on each.

1. Usage: `eqscplotLoc(x, y, ratio = 1, tol = 0.04, uin, ...)`
2. Arguments:

- `x`: vector of `x` values, or a 2-column matrix, or a list with components “`x`” and “`y`”.
- `y`: vector of `y` values.
- `ratio`: desired ratio of units on the axes. Units on the `y` axis are drawn at “`ratio`” times the size of units on the `x` axis. Ignored if “`uin`” is specified and of length 2.
- `tol`: proportion of white space at the margins of plot.
- `uin`: desired values for the units-per-inch parameter. If of length 1, the desired units per inch on the `x` axis.
- `...`: further arguments for “`plot`” and graphical parameters. Note that “`par(xaxs="i", yaxs="i")`” is enforced, and “`xlim`” and “`ylim`” will be adjusted accordingly.

### 3. Return Values:

- invisibly, the values of “`uin`” used for the plot.

Plots with Geometrically Equal Scales `eqsplot` from the package `MASS` (Version: 7.3-3, 2009-10-15 Maintainer: Brian Ripley <ripley@stats.ox.ac.uk>).

To make the package stand alone, this function is included as source.

Limits for the `x` and `y` axes are chosen so that they include the data. One of the sets of limits is then stretched from the midpoint to make the units in the ratio given by “`ratio`”. Finally both are stretched by “`1 + tol`” to move points away from the axes, and the points plotted.

Implementation in R .

### REFERENCE:

- Venables, W. N. and Ripley, B. D., (2002), “Modern Applied Statistics with S”, Fourth edition, Springer.

## F Data Sets

### F.1 Breast\_A

Microarray data set of van’t Veer breast cancer.

Microarray data from Broad Institute “Cancer Program Data Sets” which was produced by (van’t Veer *et al.*, 2002) (<http://www.broadinstitute.org/cgi-bin/cancer/datasets.cgi>) Array S54 was removed because it is an outlier.

Goal was to find a gene signature to predict the outcome of a cancer therapy that is to predict whether metastasis will occur. A 70 gene signature has been discovered.

Here we want to find subclasses in the data set.

(Hoshida *et al.*, 2007) found 3 subclasses and verified that 50/61 cases from class 1 and 2 were ER positive and only in 3/36 from class 3.

XBreast is the data set with 97 samples and 1213 genes, CBreast give the three subclasses from (Hoshida *et al.*, 2007).



## F.2 DLBCL\_B

Microarray data set of Rosenwald diffuse large-B-cell lymphoma.

Microarray data from Broad Institute “Cancer Program Data Sets” which was produced by (Rosenwald *et al.*, 2002) (<http://www.broadinstitute.org/cgi-bin/cancer/datasets.cgi>)

Goal was to predict the survival after chemotherapy

(Hoshida *et al.*, 2007) divided the data set into three classes:

- OxPhos: oxidative phosphorylation
- BCR: B-cell response
- HR: host response

We want to identify these subclasses.

The data has 180 samples and 661 probe sets (genes).

DLBCL is the data set with 180 samples and 661 genes, DLBCL give the three subclasses from (Hoshida *et al.*, 2007).

## F.3 Multi\_A

Microarray data set of Su on different mammalian tissue types.

Microarray data from Broad Institute “Cancer Program Data Sets” which was produced by (Su *et al.*, 2002) (<http://www.broadinstitute.org/cgi-bin/cancer/datasets.cgi>)

Gene expression from human and mouse samples across a diverse array of tissues, organs, and cell lines have been profiled. The goal was to have a reference for the normal mammalian transcriptome.

Here we want to identify the subclasses which correspond to the tissue types.

The data has 102 samples and 5565 probe sets (genes).

Multi is the data set with 102 samples and 5565 genes, Multi give the four subclasses corresponding to the tissue types.

## References

- Bithas, P. S., Sagias, N. C., Tsiftsis, T. A., and Karagiannidis, G. K. (2007). Distributions involving correlated generalized gamma variables. In *Proc. Int. Conf. on Applied Stochastic Models and Data Analysis*, volume 12, Chania.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum Likelihood from Incomplete Data via the EM Algorithm. *J. R. Stat. Soc. B Met.*, **39**(1), 1–22.
- Girolami, M. (2001). A variational method for learning sparse and overcomplete representations. *Neural Comput.*, **13**(11), 2517–2532.
- Hochreiter, S., Clevert, D.-A., and Obermayer, K. (2006). A new summarization method for Affymetrix probe level data. *Bioinformatics*, **22**(8), 943–949.

- Hochreiter, S., Bodenhofer, U., Heusel, M., Mayr, A., Mitterecker, A., Kasim, A., Khamiakova, T., Sanden, S. V., Lin, D., Talloen, W., Bijmens, L., Göhlmann, H. W. H., Shkedy, Z., and Clevert, D.-A. (2010). FABIA: Factor analysis for bicluster acquisition. *Bioinformatics*. doi:10.1093/bioinformatics/btq227, <http://bioinformatics.oxfordjournals.org/cgi/content/abstract/btq227>.
- Hoshida, Y., Brunet, J.-P., Tamayo, P., Golub, T. R., and Mesirov, J. P. (2007). Subclass mapping: Identifying common subtypes in independent disease data sets. *PLoS ONE*, **2**(11), e1195.
- Hoyer, P. O. (2004). Non-negative matrix factorization with sparseness constraints. *J. Mach. Learn. Res.*, **5**, 1457–1469.
- Hyvärinen, A. and Oja, E. (1999). A fast fixed-point algorithm for independent component analysis. *Neural Comput.*, **9**(7), 1483–1492.
- Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, **401**(6755), 788–791.
- Lee, D. D. and Seung, H. S. (2001). Algorithms for non-negative matrix factorization. In *Advances in Neural Information Processing Systems 13*, pages 556–562.
- Paatero, P. and Tapper, U. (1997). Least squares formulation of robust non-negative factor analysis. *Chemometr. Intell. Lab.*, **37**, 23–35.
- Palmer, J., Wipf, D., Kreutz-Delgado, K., and Rao, B. (2006). Variational EM algorithms for non-Gaussian latent variable models. In *Advances in Neural Information Processing Systems 18*, pages 1059–1066.
- Rosenwald, A., Wright, G., Chan, W. C., Connors, J. M., Campo, E., Fisher, R. I., Gascoyne, R. D., Muller-Hermelink, H. K., Smeland, E. B., Giltmane, J. M., Hurt, E. M., Zhao, H., Averett, L., Yang, L., Wilson, W. H., Jaffe, E. S., Simon, R., Klausner, R. D., Powell, J., Duffey, P. L., Longo, D. L., Greiner, T. C., Weisenburger, D. D., Sanger, W. G., Dave, B. J., Lynch, J. C., Vose, J., Armitage, J. O., Montserrat, E., L'opez-Guillermo, A., Grogan, T. M., Miller, T. P., LeBlanc, M., Ott, G., Kvaloy, S., Delabie, J., Holte, H., Krajci, P., Stokke, T., and Staudt, L. M. (2002). The use of molecular profiling to predict survival after chemotherapy for diffuse large-B-cell lymphoma. *New Engl. J. Med.*, **346**, 1937–1947.
- Samperi, D. (2006). Building R packages that call C++ functions  
Rcpp: R/C++ interface classes version 5.0. Vignette of a CRAN package. <http://cran.r-project.org/web/packages/Rcpp/vignettes/RcppAPI.pdf>.
- Su, A. I., Cooke, M. P., A.Ching, K., Hakak, Y., Walker, J. R., Wiltshire, T., Orth, A. P., Vega, R. G., Sapinoso, L. M., Moqrich, A., Patapoutian, A., Hampton, G. M., Schultz, P. G., and Hogenesch, J. B. (2002). Large-scale analysis of the human and mouse transcriptomes. *P. Natl. Acad. Sci. USA*, **99**(7), 4465–4470.
- Talloen, W., Clevert, D.-A., Hochreiter, S., Amaratunga, D., Bijmens, L., Kass, S., and Göhlmann, H. W. H. (2007). I/NI-calls for the exclusion of non-informative genes: a highly effective feature filtering tool for microarray data. *Bioinformatics*, **23**(21), 2897–2902.
- van't Veer, L. J., Dai, H., van de Vijver, M. J., He, Y. D., Hart, A. A., Mao, M., Peterse, H. L., van der Kooy, K., Marton, M. J., Witteveen, A. T., Schreiber, G. J., Kerkhoven, R. M., Roberts, C., Linsley, P. S., Bernards, R., and Friend, S. H. (2002). Gene expression profiling predicts clinical outcome of breast cancer. *Nature*, **415**, 530–536.