

Using Package NMF

Renaud Gaujoux, <renaud@cbio.uct.ac.za>

March 15, 2010

This vignette presents the **NMF** package, which implements a framework for Nonnegative Matrix Factorization (NMF) algorithms in R [R Software, 2008]. The objective is to provide an implementation of some standard algorithms, while allowing the user to easily implement new methods that integrate into the package's framework.

Contents

1	Overview	2
1.1	Package features	2
1.2	Nonnegative Matrix Factorization	2
1.3	Algorithms	3
1.4	Initialization: seeding methods	3
1.5	How to run NMF algorithms	4
2	Use case: Golub dataset	5
2.1	Single run	5
2.1.1	Performing a single run	5
2.1.2	Handling the result	5
2.1.3	Selecting the features	8
2.2	Specifying the algorithm	8
2.2.1	Custom function	10
2.3	Specifying the seeding method	10
2.3.1	Named seeding method	10
2.3.2	Numerical seed	11
2.3.3	Fixed factorization	11
2.3.4	Custom function	12
2.4	Multiple runs	12
2.4.1	Sequential runs	12
2.4.2	Parallel computing on multi-core machines	13
2.4.3	High Performance Computing on a cluster	14
2.5	Estimating the factorization rank	15
2.6	Visualization methods	17
2.7	Comparing algorithms	19
3	Advanced usage	21
3.1	Custom algorithm	22
3.1.1	Using a custom algorithm	22
3.1.2	Using a custom distance measure	24
3.1.3	Handling mixed sign data	25
3.1.4	Specifying the NMF model	26
3.2	Custom seeding method	27
3.3	Package specific options	28

The last stable version of the NMF package can be installed from any [CRAN](#) repository mirror, , and loaded with the standard calls:

```
## Not run
# install.packages('NMF')
library(NMF)
```

1 Overview

1.1 Package features

This section provides a quick overview of the NMF package's features. Section 2 provides more details, as well as sample code on how to actually perform common tasks in NMF analysis.

The NMF package:

- 8 built-in algorithms;
- 4 built-in seeding methods;
- Single interface to perform all algorithms, and combine them with the seeding methods;
- Provides a common framework to test, compare and develop NMF methods;
- Accept custom algorithms and seeding methods;
- Plotting utility functions to visualize and help in the interpretation of the results;
- Transparent parallel computations;
- Optional layer for bioinformatics based on BioConductor [3];

1.2 Nonnegative Matrix Factorization

This section gives a formal definition for Nonnegative Matrix Factorization problems, and defines the notations used throughout the vignette.

Let X be a $n \times p$ non-negative matrix, (i.e with $x_{ij} \geq 0$, denoted $X \geq 0$), and $r > 0$ an integer. Non-negative Matrix Factorization (NMF) consists in finding an approximation

$$X \approx WH, \quad (1)$$

where W, H are $n \times r$ and $r \times p$ non-negative matrices, respectively. In practice, the factorization rank r is often chosen such that $r \ll \min(n, p)$. The objective behind this choice is to summarize and split the information contained in X into r factors: the columns of W .

Depending on the application field, these factors are given different names: basis images, metagenes, source signals. In this vignette we equivalently and alternatively use the terms *basis matrix* or *metagenes* to refer to matrix W , and *mixture coefficient matrix* and *metagene expression profiles* to refer to matrix H .

The main approach to NMF is to estimate matrices W and H as a local minimum:

$$\min_{W, H \geq 0} \underbrace{[D(X, WH) + R(W, H)]}_{=F(W, H)} \quad (2)$$

where

- D is a loss function that measures the quality of the approximation. Common loss functions are based on either the Frobenius distance

$$D : A, B \mapsto \text{Tr}(AB^t) = \frac{1}{2} \sum_{ij} (a_{ij} - b_{ij})^2,$$

or the Kullback-Leibler divergence.

$$D : A, B \mapsto \sum_{i,j} a_{ij} \log \frac{a_{ij}}{b_{ij}} - a_{ij} + b_{ij}.$$

- R is an optional regularization function, defined to enforce desirable properties on matrices W and H , such as smoothness or sparsity [A. Cichocki *et al.*, 2004].

1.3 Algorithms

NMF algorithms generally solve problem (2) iteratively, by building a sequence of matrices (W_k, H_k) that reduces at each step the value of the objective function F . Beside some variations in the specification of F , they also differ in the optimization techniques that are used to compute the updates for (W_k, H_k) .

For reviews on NMF algorithms see [Berry *et al.*, 2006, Chu *et al.*, 2004] and references therein.

The NMF package implements a number of published algorithms, and provides a general framework to implement other ones.

The built-in algorithms are listed or retrieved with function `nmfAlgorithm`. A given algorithm is retrieved by its name (a `character` key), that is partially matched against the list of available algorithms:

```
# list all available algorithms
nmfAlgorithm()
      [1] "brunet" "lee"    "lnmf"   "nsNMF"  "offset" "pe-nmf" "snmf/l" "snmf/r"

# retrieve a specific algorithm: 'brunet'
nmfAlgorithm('brunet')
      <object of class: NMFStrategyIterative >
      name:          brunet
      objective:      'KL'
      NMF model:      NMFstd
      <Iterative schema:>
      Update : 'nmf.update.brunet'
      Stop : 'nmf.stop.consensus'
      WrapNMF : ''

# partial match is also fine
identical(nmfAlgorithm('br'), nmfAlgorithm('brunet'))
      [1] TRUE
```

1.4 Initialization: seeding methods

NMF algorithms need to be initialized with a seed (i.e. a value for W_0 and/or H_0 ¹), from which to start the iteration process. Because there is no global minimization algorithm, and due to

¹Some algorithms only need one matrix factor (either W or H) to be initialized. See for example the SNMF/R(L) algorithm of Kim and Park [Kim and Park, 2007].

the problem's high dimensionality, the choice of the initialization is in fact very important to ensure meaningful results.

The more common seeding method is to use a random starting point, where the entries of W and/or H are drawn from a uniform distribution, usually within the same range as the target matrix's entries. This method is very simple to implement. However, a major drawback is that to achieve stability it requires to perform multiple runs, each with a different starting point. This significantly increases the computation time needed to obtain the desired factorization.

To tackle this problem, some methods have been proposed so as to compute a reasonable starting point from the target matrix itself. The objective is to produce deterministic algorithms that need to run only once, still giving meaningful results.

For a review on some existing NMF initializations see [Albright et al., 2006] and references therein.

The `NMF` package implements a number of already published seeding methods, and provides a general framework to implement other ones.

The built-in seeding methods are listed or retrieved with function `nmfSeed`. A given seeding method is retrieved by its name (a `character` key) that is partially matched against the list of available seeding methods:

```
# list all available seeding methods
nmfSeed()
[1] "ica"      "nndsvd" "none"    "random"

# retrieve a specific method: 'nndsvd'
nmfSeed('nndsvd')
<object of class: NMFSeed >
name:      nndsvd
method:    <function>

# partial match is also fine
identical(nmfSeed('nn'), nmfSeed('nndsvd'))
[1] TRUE
```

1.5 How to run NMF algorithms

Method `nmf` provides a single interface to run NMF algorithms. It can directly perform NMF on object of class `matrix` or `data.frame` and `ExpressionSet` – if the `Biobase` package is installed. The interface has four main parameters:

```
nmf(x, rank, method, seed, ...)
```

`x` is the target `matrix`, `data.frame` or `ExpressionSet` ²

`rank` is the factorization rank, i.e. the number of columns in matrix W .

`method` is the algorithm used to estimate the factorization. The default algorithm is given by the package specific option `'default.algorithm'`, which defaults to `'brunet'` on installation [Brunet et al., 2004].

`seed` is the seeding method used to compute the starting point. The default method is given by the package specific option `'default.seed'`, which defaults to `'random'` on initialization (see method `?nmf` for details on its implementation).

See also `?nmf` for details on the interface and extra parameters.

²`ExpressionSet` is the base class for handling microarray data in BioConductor, and is defined in the `Biobase` package.

2 Use case: Golub dataset

We illustrate the functionalities and the usage of the NMF package, by analysing the – now standard – Golub dataset on leukemia. It was used in several papers on NMF [Brunet *et al.*, 2004, 2] and is included in the NMF package’s data, wrapped into an `ExpressionSet` object. For performance reason we use here only the first 200 genes:

```
data(esGolub)
esGolub
  ExpressionSet (storageMode: lockedEnvironment)
  assayData: 5000 features, 38 samples
    element names: exprs
  phenoData
    sampleNames: ALL_19769_B-cell, ALL_23953_B-cell, ..., AML_7 (38 total)
    varLabels and varMetadata description:
      Sample: Sample name from the file ALL_AML_data.txt
      ALL.AML: ALL/AML status
      Cell: Cell type
  featureData
    featureNames: M12759_at, U46006_s_at, ..., D86976_at (5000 total)
    fvarLabels and fvarMetadata description:
      Description: Short description of the gene
  experimentData: use 'experimentData(object)'
  Annotation:

esGolub <- esGolub[1:50,]
```

Note: To run this example, the `Biobase` package from BioConductor is required.

2.1 Single run

2.1.1 Performing a single run

To run the default NMF algorithm on data `esGolub` with a factorization rank of 3, we call:

```
# default NMF algorithm
res <- nmf(esGolub, 3)
```

Here we did not specify either the algorithm or the seeding method, so that the computation is done using the default algorithm and is seeded by the default seeding methods. These defaults are set in the package specific options `'default.algorithm'` and `'default.seed'` respectively.

See also sections 2.2 and 2.3 for how to explicitly specify the algorithm and/or the seeding method.

2.1.2 Handling the result

The result of a single NMF run is an object of class `NMFfit`, that holds both the fitted NMF model and data about the run:

```
res
```

```

<Object of class: NMFfit >
# Model:
  <Object of class: NMFstd >
    features: 50
    basis/rank: 3
    samples: 38
# Details:
  algorithm: brunet
  seed: random
  distance metric: 'KL'
  residuals: 178467.1
  Iterations: 640
Timing:
  user  system elapsed
  0.690   0.000   0.684

```

The fitted model can be retrieved via method `fit`, which returns an object of class `NMF`:

```

fit(res)
  <Object of class: NMFstd >
    features: 50
    basis/rank: 3
    samples: 38

```

The estimated target matrix can be retrieved via the generic method `fitted`, which returns a – generally big – matrix:

```

V.hat <- fitted(res)
dim(V.hat)
[1] 50 38

```

Quality and performance measures about the factorization are computed by method `summary`:

```

summary(res)
      rank sparseness.basis sparseness.coef      niter
3.000000e+00  6.829776e-01  5.775976e-01  6.400000e+02
      time      residuals
6.900000e-01  1.784671e+05

# More quality measures are computed, if the target matrix is provided:
summary(res, target=esGolub)
      rank sparseness.basis sparseness.coef      rss
3.000000e+00  6.829776e-01  5.775976e-01  5.016224e+08
      evar      niter      time      residuals
8.631025e-01  6.400000e+02  6.900000e-01  1.784671e+05

```

If there is some prior knowledge of classes present in the data, some other measures about the unsupervised clustering's performance are computed (purity, entropy, ...). Here we use the

phenotypic variable `Cell` found in the Golub dataset, that gives the samples' cell-types (it is a factor with levels: T-cell, B-cell or NA):

```
summary(res, class=esGolub$Cell)
```

	rank	sparseness.basis	sparseness.coef	purity
	3.000000e+00	6.829776e-01	5.775976e-01	7.105263e-01
	entropy	niter	time	residuals
	5.849030e-01	6.400000e+02	6.900000e-01	1.784671e+05

The basis matrix (i.e. matrix W or the metagenes) and the mixture coefficient matrix (i.e matrix H or the metagene expression profiles) are retrieved using methods `basis` and `coef` respectively:

```
# get matrix W
w <- basis(res)
dim(w)
[1] 50 3

# get matrix H
h <- coef(res)
dim(h)
[1] 3 38
```

If one wants to keep only part of the factorization, one can directly subset on the NMF object on features and samples (separately or simultaneously):

```
# keep only the first 100 features TODO: put that to 100
res[1:100,]
<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 100
basis/rank: 3
samples: 38
# Details:
algorithm: brunet
seed: random
distance metric: 'KL'
residuals: 178467.1
Iterations: 640
Timing:
  user system elapsed
  0.690 0.000 0.684

# keep only the first 10 samples
res[,1:10]
<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 50
basis/rank: 3
```

```

samples: 10
# Details:
algorithm: brunet
seed: random
distance metric: 'KL'
residuals: 178467.1
Iterations: 640
Timing:
  user  system elapsed
 0.690   0.000   0.684

```

```
[1] 20 10 3
```

2.1.3 Selecting the features

In general NMF matrix factors are sparse, so that the metagenes can usually be characterized by a relatively small set of genes. Those are determined based on their relative contribution to each metagene.

Kim and Park [Kim and Park, 2007] defined a procedure to extract the relevant genes for each metagene, based on a gene scoring schema.

The NMF package implements this procedure in methods `featureScore` and `extractFeature`:

```

# only compute the scores
s <- featureScore(res)
summary(s)
      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
2.995e-05 3.849e-02 1.185e-01 2.130e-01 2.829e-01 9.997e-01

# compute the scores and characterize each metagene
s <- extractFeatures(res)
str(s)
List of 3
 $ 1: int 30
 $ 2: int [1:2] 2 39
 $ 3: int [1:2] 1 43
 - attr(*, "threshold")= num 0.505

```

2.2 Specifying the algorithm

subsubsectionNamed algorithm The NMF package provides a number of built-in algorithms, that are listed or retrieved by function `nmfAlgorithm`. Each algorithm is identified by a unique name. The following algorithms are currently implemented:

```

nmfAlgorithm()
[1] "brunet" "lee"    "lnmf"   "nsNMF"  "offset" "pe-nmf" "snmf/l" "snmf/r"

```


The algorithm used to compute the NMF is specified in the third argument (`method`). For example, to use the Lee and Seung [Lee and Seung, 2000] NMF algorithm based on the Frobenius euclidean norm, one make the following call:

```
# using Lee and Seung's algorithm
res <- nmf(esGolub, 3, 'lee')
res
<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 50
basis/rank: 3
samples: 38
# Details:
algorithm: lee
seed: random
distance metric: 'euclidean'
residuals: 183823790
Iterations: 570
Timing:
  user system elapsed
  0.770  0.000  0.763
```

To use the Nonsmooth NMF algorithm from [Pascual-Montano *et al.*, 2006]:

```
# using the Nonsmooth NMF algorithm with parameter theta=0.7
res <- nmf(esGolub, 3, 'ns', theta=0.7)
res
<Object of class: NMFfit >
# Model:
<Object of class: NMFns >
features: 50
basis/rank: 3
samples: 38
theta: 0.7
# Details:
algorithm: nsNMF
seed: random
distance metric: 'KL'
residuals: 224297.2
Iterations: 820
Timing:
  user system elapsed
  1.090  0.000  1.116
```

Or to use the PE-NMF algorithm from [5]:

```
# using the PE-NMF algorithm with parameters alpha=0.01, beta=1
res <- nmf(esGolub, 3, 'pe', alpha=0.01, beta=1)
res
```

```

<Object of class: NMFfit >
# Model:
  <Object of class: NMFstd >
    features: 50
    basis/rank: 3
    samples: 38
# Details:
  algorithm: pe-nmf
  seed: random
  distance metric: <function>
  residuals: 52.44378
  parameters:
    $alpha
    [1] 0.01

    $beta
    [1] 1

  Iterations: 2000
  Timing:
    user  system elapsed
    1.400   0.000   1.392

```

2.2.1 Custom function

The NMF package provides the user the possibility to define his own algorithm, and benefit from all the functionalities available in the NMF framework. There are only few constraints on the way the custom algorithm must be defined. See the details in Section [3.1.1](#).

2.3 Specifying the seeding method

The seeding method used to compute the starting point for the chosen algorithm can be set via argument `seed`. Note that if the seeding method is deterministic there is no need to perform multiple run anymore.

2.3.1 Named seeding method

Similarly to the algorithms, the `nmfSeed` function can be used to list or retrieve the built-in seeding methods.

The following seeding methods are currently implemented:

```

nmfSeed()
[1] "ica"      "nndsvd"  "none"    "random"

```

To use a specific method to seed the computation of a factorization, one can provide the name of the seeding method:

```

res <- nmf(esGolub, 3, seed='nndsvd')
res
<Object of class: NMFfit >
# Model:

```

```

<Object of class: NMFstd >
features: 50
basis/rank: 3
samples: 38
# Details:
algorithm: brunet
seed: nndsvd
distance metric: 'KL'
residuals: 173800.5
Iterations: 1130
Timing:
  user system elapsed
  1.110   0.000   1.106

```

2.3.2 Numerical seed

Another possibility, useful when comparing methods or testing the reproducibility of the results, is to set the seed of the random generator by passing a numerical value in argument `seed`. This will call the function `set.seed` from package `base` before using the `'random'` seeding method:

```

res <- nmf(esGolub, 3, seed=123456)
res
<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 50
basis/rank: 3
samples: 38
# Details:
algorithm: brunet
seed: 123456
distance metric: 'KL'
residuals: 178467.1
Iterations: 640
Timing:
  user system elapsed
  0.63   0.00   0.63

```

2.3.3 Fixed factorization

Yet another option is to completely specify the initial factorization, by passing values for matrices W and H :

```

n <- nrow(esGolub); p <- ncol(esGolub)
res <- nmf(esGolub, 3, seed=NULL, W=matrix(0.5, n, 3), H=matrix(0.3, 3, p))
res
<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 50
basis/rank: 3

```

```

samples: 38
# Details:
algorithm: brunet
seed: none
distance metric: 'KL'
residuals: 300070.8
Iterations: 420
Timing:
  user  system elapsed
0.410   0.000   0.408

```

Important: in this case, argument `seed` must absolutely be set to `NULL`, otherwise the model instantiated with matrices W and H would only be used as a template, and reset passing it to the default seeding method.

Two alternative ways of doing this would be to pass matrices W and H through argument `model`, or a NMF model to argument `seed`:

```

res <- nmf(esGolub, 3, seed=NULL
           , model=list(W=matrix(0.5, n, 3), H=matrix(0.3, 3, p)))
# or
res <- nmf(esGolub, 3, seed=newNMF(W=matrix(0.5, n, 3), H=matrix(0.3, 3, p)))

```

2.3.4 Custom function

The NMF package provides the user the possibility to define his own seeding method, and benefit from all the functionalities available in the NMF framework. There are only few constraints on the way the custom seeding method must be defined. See the details in Section~3.2.

2.4 Multiple runs

When the seeding method is stochastic, multiple runs are usually required to achieve stability or a resonable result. This can be done by setting argument `nrun` to the desired value. For performance reason we use `nrun=5` here, but a typical choice would lies between 100 and 200:

2.4.1 Sequential runs

```

res.multirun <- nmf(esGolub, 3, nrun=5)
res.multirun
  <Object of class: NMFSet >
method: brunet
runs: 5
fits: 1
Timing:
  user  system elapsed
2.140   2.130   2.531
Avg. timing:
  user  system elapsed
0.4280  0.4260  0.5062

```

As we can see from the results above, the returned object contains only one fit, from the 5 runs that were performed. Indeed the default behaviour is to only keep the factorization that achieves the lowest approximation error (i.e. the lowest objective value). Even during the computation, only the current best factorization is kept. This limits the memory requirement of performing multiple runs, which in turn allows to perform more runs.

However if one is interested in keeping the results from all the runs, one can set the option `keep.all=TRUE`:

```
# explicitly setting the option
nmf(esGolub, 3, nrun=5, .options=list(keep.all=TRUE))
# or using letter code 'k' in argument .options
nmf(esGolub, 3, nrun=5, .options='k')
```

Note that keeping all the results may be memory consuming. For example, a 3-rank NMF model³ for the Golub gene expression matrix (5000×38) takes a bit less than 350Kb.

2.4.2 Parallel computing on multi-core machines

To speed-up the analysis whenever possible, the NMF package implements transparent parallel computations when run on multi-core machines. It uses the `foreach` framework developed by R Evolution Computing [foreach, 2009], together with the related `doMC` parallel backend [doMC, 2009] – based on the `multicore` package – to make use of all the CPUs available on the system. Each core will simultaneously perform part of the runs. Therefore, the required memory increases linearly with the number of cores used. When only the best run is of interest, the memory usage is optimized by using shared memory and mutex objects from the `bigmemory` package, to only keep the current best factorization.

IMPORTANT NOTE: because it uses the `multicore` package, parallel computation over multi-cores is available only for Unix and Mac machines.

The default parallel backend used by the `nmf` function is defined by the package specific option `'parallel.backend'`, which defaults to `'mc'` – for `doMC`. The backend can also be set on runtime via argument `'.pbackend'`.

There are two other runtime options, `parallel` and `parallel.required`, that can be passed via argument `.options`, to control the behaviour of the parallel computation (see below).

A call for multiple runs will be computed in parallel if one of the following condition is satisfied:

- call with option `'P'` or `parallel.required` set to `TRUE` (note the upper case in `'P'`). In this case, if for any reason the computation cannot be run in parallel (packages requirements, OS, ...), then an error is thrown. Use this mode to force the parallel execution.
- call with option `'p'` or `parallel` set to `TRUE`. In this case if something prevents a parallel computation, the factorizations will be done sequentially.
- a valid parallel backend is specified in argument `.pbackend`. For the moment can either be the string `'mc'` or a single `numeric` value specifying the number of core to use. Unless option `'P'` is specified, it will run using option `'p'` (i.e. try-parallel mode).

³i.e. the result of a single NMF run with rank equal 3.

Examples

```
# default call will try to run in parallel using all the cores
# => will be in parallel if all the requirements are satisfied
nmf(esGolub, 3, nrun=10, .opt='v')
Runs: 1 2 3 4 5 6 7 8 9 10 ... DONE
<Object of class: NMFSets>
method: brunet
runs: 10
fits: 1
Timing:
  user  system elapsed
6.070   6.050   4.576
Avg. timing:
  user  system elapsed
0.6070  0.6050  0.4576
```

```
# specifying the number of cores to use
nmf(esGolub, 3, nrun=10, .opt='v', .pbackend=2)
# force parallel computation: use option 'P'
nmf(esGolub, 3, nrun=10, .opt='vP')
# force sequential computation: use option '-p'
nmf(esGolub, 3, nrun=10, .opt='v-p')
# or use the SEQ backend: .pbackend=NULL or 'seq'
nmf(esGolub, 3, nrun=10, .opt='vp', .backend=NULL)
```

2.4.3 High Performance Computing on a cluster

To achieve further speed-up, the computation can be run on an HPC cluster. In our tests we used the `doMPI` package to perform 100 factorizations using hybrid parallel computation on 4 quadri-core machines – i.e. making use of all the cores computation on each machine.

The scripts used to launch and run the factorizations can be found in file `mpi.R` in the package's `examples` directory:

```
file.show(file.system('examples/mpi.R', package='NMF'))
# and
file.show(file.system('examples/mpi_run.sh', package='NMF'))
```

The script file `mpi.R` contains some extra code to log and trace the computation. Reducing it to the essential gives the following piece of code:

```
## 0. Create and register an MPI cluster
library(doMPI)
cl <- startMPIcluster()
registerDoMPI(cl)
library(NMF)
```

```

## 1. Schedule the runs accross the workers
nrun <- 100;
nworker <- getDoParWorkers();
ntasks <- rep(round(nrun/nworker), nworker)
# allocate remainder runs
if( (remain <- nrun %% nworker) > 0 )
  ntasks[1:remain] <- ntasks[1:remain] + 1
## 2. Send the jobs to the workers using a foreach loop
t <- system.time({
  res <- foreach(i=1:getDoParWorkers(), n=ntasks,
    .packages = c('NMF', 'doMC', 'Biobase')) %dopar% {

    # each worker run its factorizations in parallel
    #Note: only the best result is kept
    data(esGolub)
    nmf(esGolub, 3, 'brunet', nrun=n, .opt='p')

  }
})
## 3. reduce the result and save it in a file
res <- NMF:::join(res, runtime=t, .merge=TRUE)
save(res, file='result.RData')
## 4. Shutdown the cluster and quit MPI
closeCluster(cl)
mpi.quit()

```

Passing the following shell script to *qsub* should launch the execution on a Sun Grid Engine HPC cluster, with OpenMPI. Some adaptation might be necessary for other queueing systems.

```

#!/bin/bash
#$ -cwd
#$ -q opteron.q
#$ -pe mpich 5
echo "Got $NSLOTS slots. $TMP/machines"

orterun -v -n $NSLOTS -hostfile $TMP/machines R --slave -f mpi.R

```

2.5 Estimating the factorization rank

A critical parameter in NMF is the factorization rank r . It defines the number of metagenes used to approximate the target matrix. Given a NMF method and the target matrix, a common way of deciding on r is to try different values, compute some quality measure of the results, and choose the best value according to this quality criteria.

The NMF package provides functions to run this procedure and plot the quality measures. Note that this can be lengthy. Whereas the standard NMF procedure usually involves several hundreds of random initialization, performing 30-50 runs is considered sufficient to get a robust estimate of the factorization rank [Brunet *et al.*, 2004, 4]. For performance reason, we perform here only 10 runs for each value of the rank.

```

# perform 10 runs for each value of r in range 2:6
estim.r <- nmfEstimateRank(esGolub, range=2:6, nrun=10, seed=123456)

```

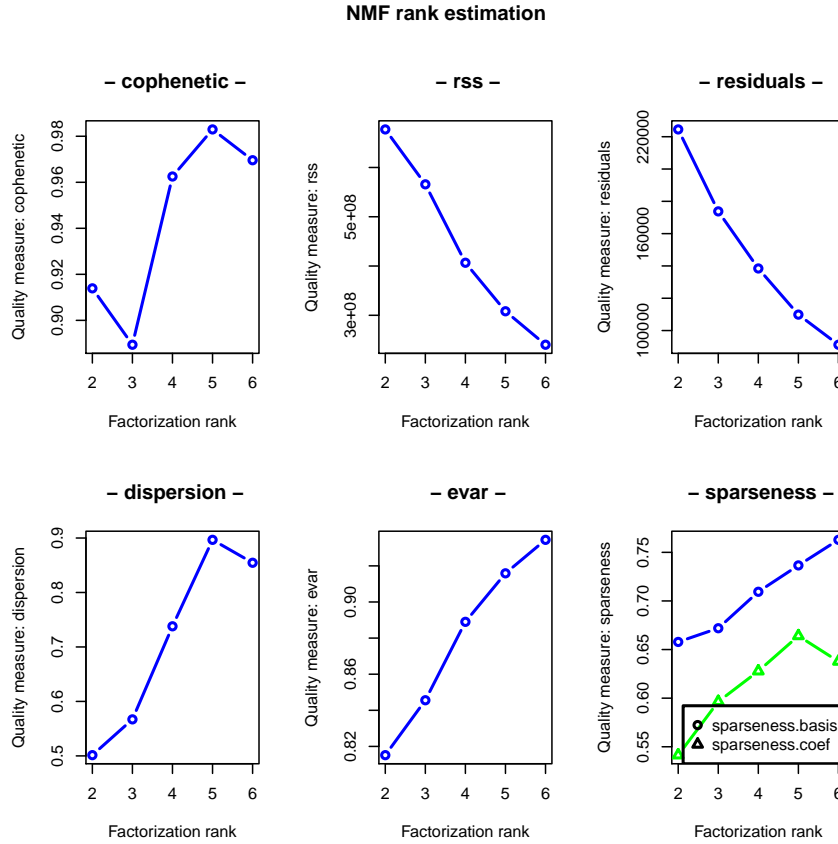


Figure 1: Estimation of the rank: Quality measures computed from 10 runs for each value of r

The result is a S3 object of class `NMF.rank`, that contains a `data.frame` with the quality measures in column, and the values of r in row. It also contains a list of the consensus matrix for each value of r .

All the measures can be plotted at once by the following call, the result is shown in Figure~2.5:

```
plot(estim.r)
```

Several approaches have been proposed to choose the optimal value of r . For example, [Brunet *et al.*, 2004] proposed to take the first value of r for which the cophenetic coefficient starts decreasing, [4] suggested to choose the first value where the RSS curve presents an inflection point, and [1] considered the smallest value at which the decrease in the RSS is lower than the decrease of the RSS obtained from random data.

The former approach may be useful to prevent or detect overfitting as it take into account the results for unstructured data. However it requires to compute the quality measure(s) for the random data. The NMF package provides a function that shuffles the original data, by permuting the rows of each column, using each time a different permutation. The estimation procedure can then be repeated on the randomized data, and the “random” measures can be conveniently added to the plot for comparison (see Figure~):

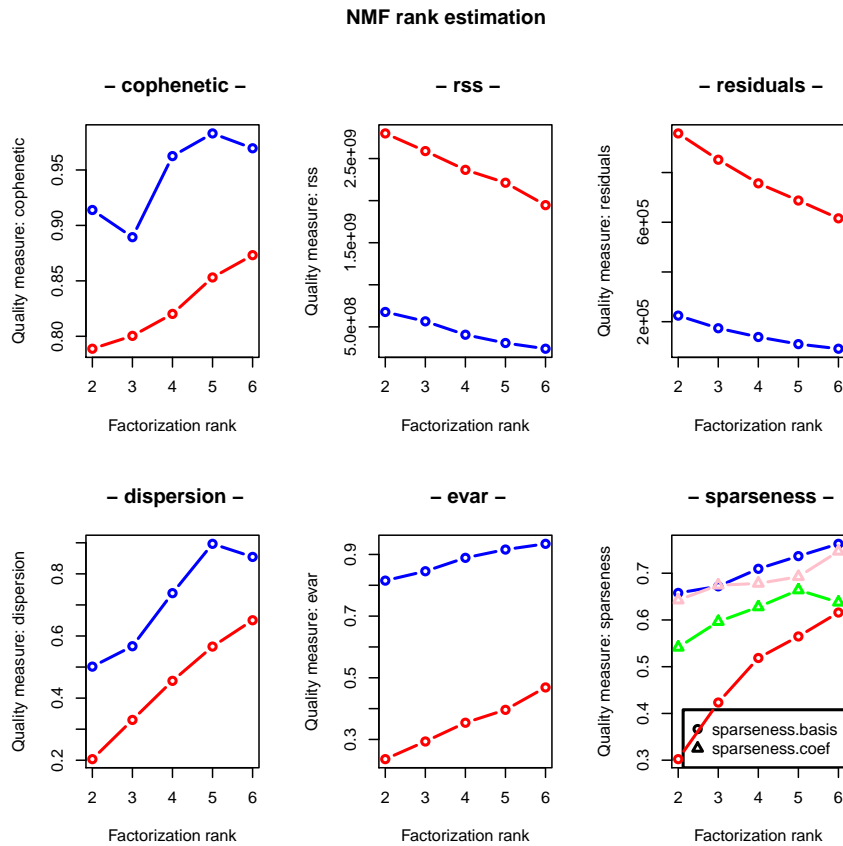


Figure 2: Estimation of the rank: Comparison of the quality measures with those obtained from randomized data

```
# shuffle original data
V.random <- randomize(esGolub)
# estimate quality measures from the shuffled data
estim.r.random <- nmfEstimateRank(V.random, range=2:6, nrun=10, seed=123456)
# plot measures on same graph
plot(estim.r, ref=estim.r.random)
```

2.6 Visualization methods

Error track

If the NMF computation is performed with error tracking enabled – using argument `.options` – the trajectory of the objective value can be plot with method `errorPlot` (see Figure 3):

```
res <- nmf(esGolub, 3, .options='t')
# or alternatively:
# res <- nmf(esGolub, 3, .options=list(track=TRUE))
errorPlot(res)
```

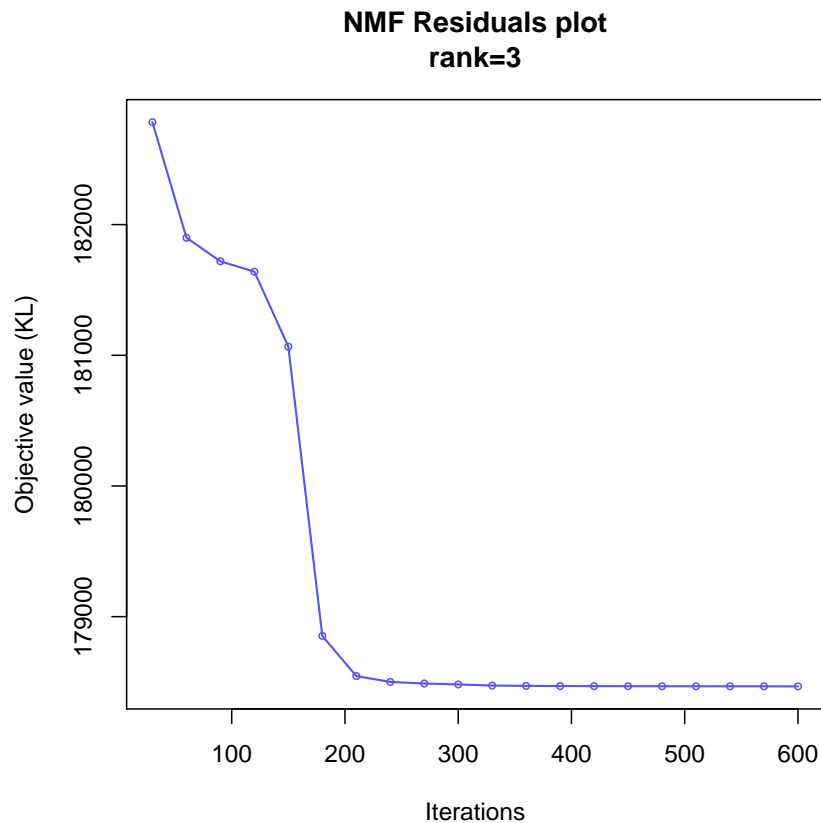


Figure 3: Error track for a single NMF run

Heatmaps

Method `metaHeatmap` provides an easy way to visualize the resulting metagenes, metaprofiles and, in the case of multiple runs, the consensus matrix. It produces pre-configured heatmaps based on function `heatmap.2` from package `gplots`. Examples of those heatmaps are shown in figures~4, 5, 6 and 7.

The following – default – call plots the metaprofiles matrix (see result Figure~4):

```
# default is to plot metaprofiles
metaHeatmap(res)
```

The metagenes matrix can be plotted specifying the second argument `what` (see result Figure~5). We use argument `filter` to select only the genes that are specific to each metagene. With `filter=TRUE`, the selection method is the one described in [Kim and Park, 2007].

```
metaHeatmap(res, what='features', filter=TRUE)
```

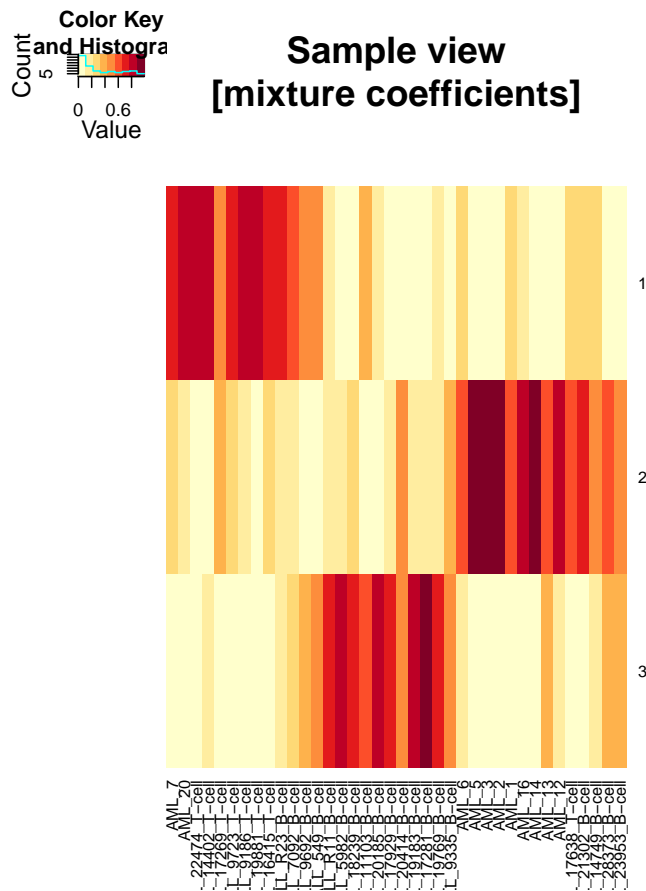


Figure 4: Heatmap of metaprofiles

In the case of multiple runs method `metaHeatmap` plots the consensus matrix, i.e. the average connectivity matrix across the runs (see results Figures 6 and 7 for a consensus matrix obtained with 100 runs of Brunet’s algorithm on Golub dataset):

```
# The cell type is used to label rows and columns
metaHeatmap(res.multirun, labRow=esGolub$Cell, labCol=esGolub$Cell)
```

2.7 Comparing algorithms

To compare the results from different algorithms, one can pass a list of methods in argument `method`. To enable a fair comparison, a deterministic seeding method should also be used. Here we fix the random seed to 123456.

```
res.multi.method <- nmf(esGolub, 3, list('brunet', 'lee', 'ns'), seed=123456)
```

Passing the result to method `compare` produces a `data.frame` that contains summary measures for each method. Again, prior knowledge of classes may be used to compute clustering quality measures:

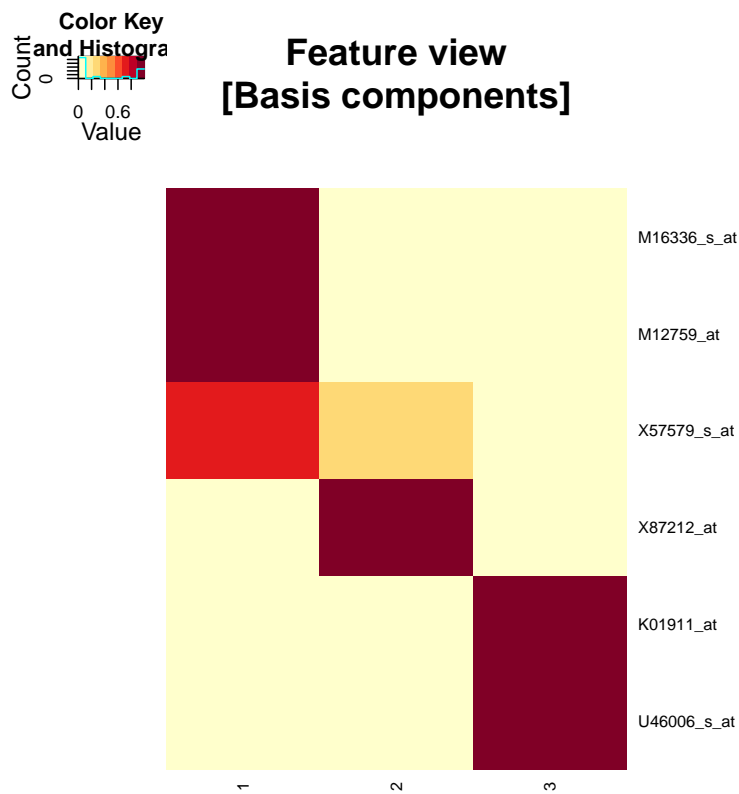


Figure 5: Heatmap of metagenes

```
compare(res.multi.method)

      method seed      metric rank sparseness.basis sparseness.coef niter
brunet brunet 123456      'KL'    3         0.6829776         0.5775976   640
nsNMF   nsNMF 123456      'KL'    3         0.7573038         0.7371265   660
lee     lee   123456 'euclidean'  3         0.7545134         0.4020160   840
      time residuals
brunet 0.67  178467.1
nsNMF  0.91  202162.5
lee    1.27 183815371.8

# If prior knowledge of classes is available
compare(res.multi.method, class=esGolub$Cell)

      method seed      metric rank sparseness.basis sparseness.coef
brunet brunet 123456      'KL'    3         0.6829776         0.5775976
nsNMF   nsNMF 123456      'KL'    3         0.7573038         0.7371265
lee     lee   123456 'euclidean'  3         0.7545134         0.4020160
      purity entropy niter time residuals
brunet 0.7105263 0.5849030 640 0.67 178467.1
nsNMF  0.6842105 0.5563767 660 0.91 202162.5
lee    0.6052632 0.7668766 840 1.27 183815371.8
```

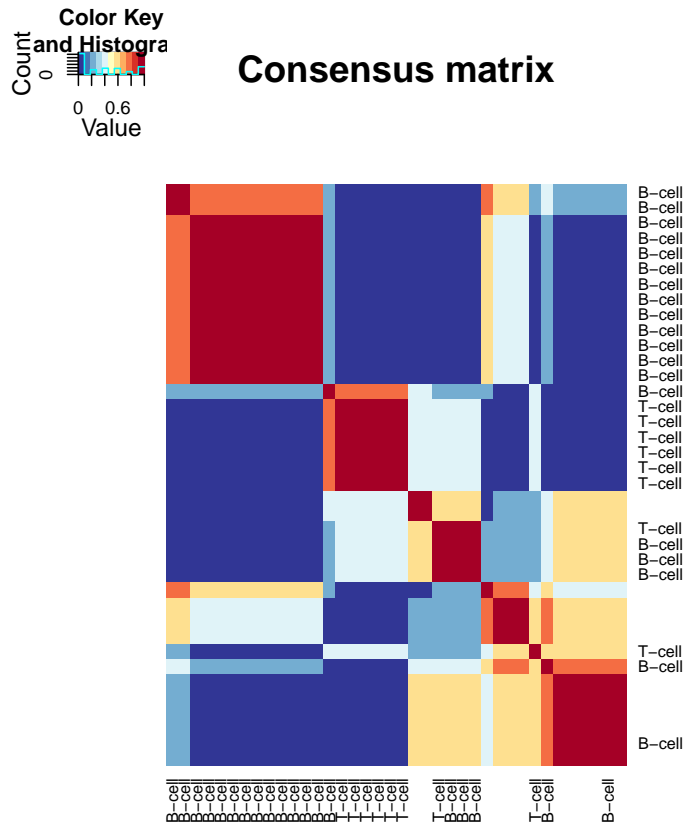


Figure 6: Heatmap of consensus matrix

When the computation is performed with error tracking enabled, an error plot is produced by method `errorplot` (see figure~8):

```
res <- nmf(esGolub, 3, list('brunet', 'lee', 'ns'), seed=123456, .options='t')
errorPlot(res)
```

3 Advanced usage

We developed the `NMF` package with the objective to facilitate the integration of new NMF methods, trying to impose only few requirements on their implementations. All the built-in algorithms and seeding methods are implemented as strategies that are called from within the main interface method `nmf`.

The user can define new strategies and those are handled in exactly the same way as the built-in ones, benefiting from the same utility functions to interpret the results and assess their performance.

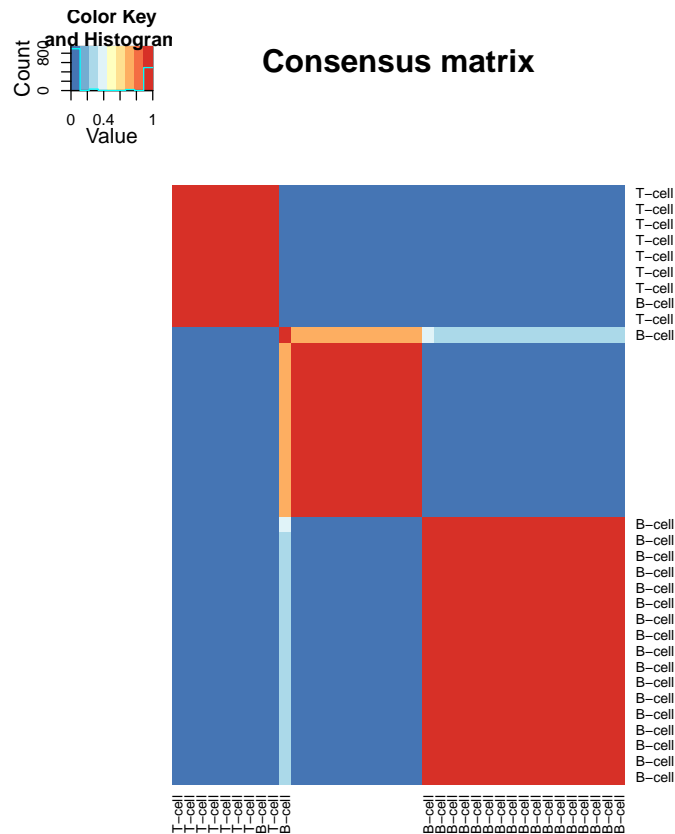


Figure 7: Heatmap of consensus matrix (100 runs of Brunet's algorithm on Golub dataset)

3.1 Custom algorithm

3.1.1 Using a custom algorithm

To define a strategy, the user needs to provide a **function** that implements the complete algorithm. It must be of the form:

```
my.algorithm <- function(x, seed, param.1, param.2){
  # do something with starting point
  # ...

  # return updated starting point
  return(seed)
}
```

Where:

target is a matrix;

start is an object that inherits from class NMF. This S4 class is used to handle NMF models (matrices W and H, objective function, etc...);

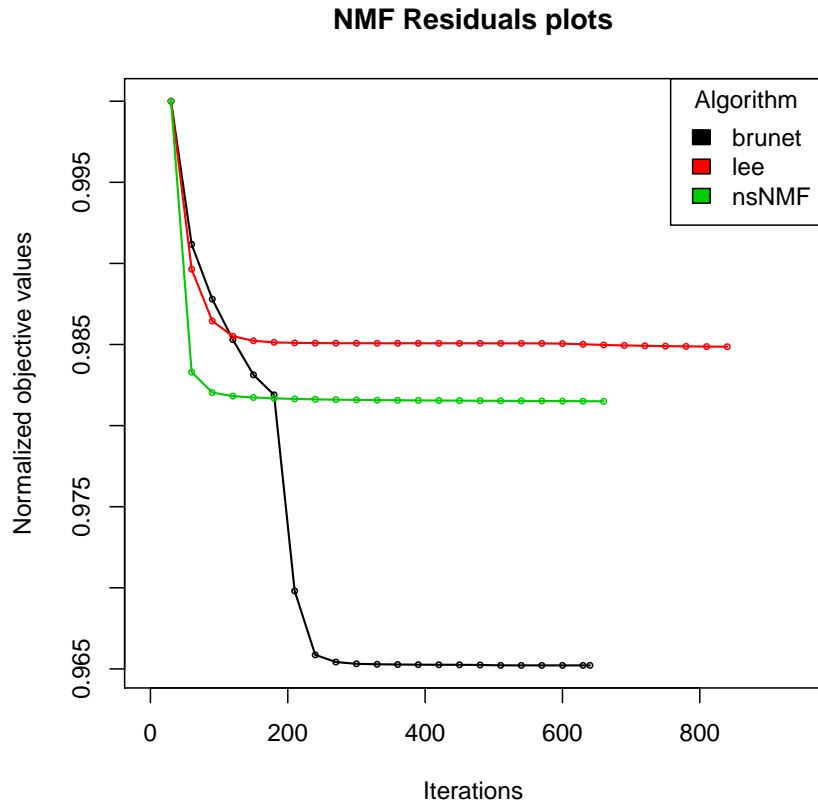


Figure 8: Error tracks comparing methods 'brunet', 'lee', 'nsNMF'

param.1, **param.2** are extra parameters specific to the algorithms;

The function must return an object that inherits from class NMF

For example:

```
my.algorithm <- function(x, seed, scale.factor=1){
  # do something with starting point
  # ...
  # for example:
  # 1. compute principal components
  pca <- prcomp(t(x), retx=TRUE)

  # 2. use the absolute values of the first PCs for the metagenes
  # Note: the factorization rank is stored in object 'start'
  factorization.rank <- nbasis(seed)
  metagenes(fit(seed)) <- abs(pca$rotation[,1:factorization.rank])
  # use the rotated matrix to get the mixture coefficient
  # use a scaling factor (just to illustrate the use of extra parameters)
  metaprofiles(fit(seed)) <- t(abs(pca$x[,1:factorization.rank])) / scale.factor

  # return updated data
```

```

    return(seed)
}

```

To use the new method within the package framework, one pass `my.algorithm` to main interface `nmf` via argument `method`. Here we apply the algorithm to some matrix `V` randomly generated:

```

n <- 50; r <- 3; p <- 20
V <-syntheticNMF(n, r, p, noise=TRUE)

```

```

nmf(V, 3, my.algorithm, scale.factor=10)
<Object of class: NMFfit >
# Model:
  <Object of class: NMFstd >
  features: 50
  basis/rank: 3
  samples: 20
# Details:
  algorithm: NMF.algo.3ca88ecf
  seed: random
  distance metric: 'euclidean'
  residuals: 698.4492
  parameters:
  $scale.factor
  [1] 10

Timing:
   user  system elapsed
 0.000   0.000   0.005

```

3.1.2 Using a custom distance measure

The default distance measure is based on the euclidean distance. If the algorithm is based on another distance measure, this one can be specified in argument `objective`, either as a `character` string corresponding to a built-in objective function, or a custom `function` definition:

```

# based on Kullbach-Leibler divergence
nmf(V, 3, my.algorithm, scale.factor=10, objective='KL')
<Object of class: NMFfit >
# Model:
  <Object of class: NMFstd >
  features: 50
  basis/rank: 3
  samples: 20
# Details:
  algorithm: NMF.algo.c058df5
  seed: random
  distance metric: 'KL'
  residuals: 1730.667

```



```

parameters:
$scale.factor
[1] 10

Timing:
  user system elapsed
0.000  0.000   0.003

<Object of class: NMFfit >
# Model:
  <Object of class: NMFstd >
features: 50
basis/rank: 3
samples: 20
# Details:
algorithm: NMF.algo.3102bbe2
seed: random
distance metric: <function>
residuals: 9.397075
parameters:
$scale.factor
[1] 10

Timing:
  user system elapsed
0.010  0.000   0.002

```

3.1.3 Handling mixed sign data

Some NMF algorithms work with relaxed constraints, where the input data and one of the matrix factors are allowed to have negative entries (eg. semi-NMF). One can plug such algorithms into the framework defined by package NMF, by specifying argument `mixed=TRUE`. The default value for argument `mixed` is `FALSE`, so that method `nmf` throws an error if applied to a matrix with some negative entries. Note that the sign of the factors is not checked, so that one can always return factors with negative entries, independently of the value of argument `mixed`. Here we reuse the previously defined custom algorithm⁴:

```

# put some negative input data
V.neg <- V; V.neg[1,] <- -1;
# this generates an error
err <- try( nmf(V.neg, 3, my.algorithm, scale.factor=10) )
err
[1] "Error in .local(x, rank, method, ...) : \n Input matrix x contains some negative entries.\n"
attr("class")
[1] "try-error"

```

⁴As it is defined here, the custom algorithm still returns nonnegative factors, which would not be desirable in a real example.

```

<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 50
basis/rank: 3
samples: 20
# Details:
algorithm: NMF.algo.3a966cd0
seed: random
distance metric: 'euclidean'
residuals: 704.3985
parameters:
$scale.factor
[1] 10

Timing:
  user system elapsed
0.010  0.000  0.003

```

3.1.4 Specifying the NMF model

If not specified in the call, the NMF model that is used is the standard one, as defined in equation (1). However, some NMF algorithms have different underlying models, such as non-smooth NMF [Pascual-Montano *et al.*, 2006] which uses an extra matrix factor that introduces an extra parameter, and change the way the target matrix is approximated.

The NMF models are defined as S4 classes that extends class `NMF`. All the available models can be retrieved using the `nmf.models()` function:

```

nmf.models()
[1] "NMFstd"      "NMFOffset"  "NMFns"

```

One can specify the NMF model to use with a custom algorithm, using argument `model`. Here we first adapt a bit the custom algorithm, to justify and illustrate the use of a different model. We use model `NMFOffset` [Badea L., 2008], that includes an offset to take into account genes that have constant expression levels across the samples:

```

my.algorithm.offset <- function(x, seed, scale.factor=1){
  # do something with starting point
  # ...
  # for example:
  # 1. compute principal components
  pca <- prcomp(t(x), retx=TRUE)

  # retrieve the model being estimated
  data.model <- fit(seed)

  # 2. use the absolute values of the first PCs for the metagenes
  # Note: the factorization rank is stored in object 'start'

```

```

    factorization.rank <- nbasis(data.model)
    metagenes(data.model) <- abs(pca$rotation[,1:factorization.rank])
    # use the rotated matrix to get the mixture coefficient
    # use a scaling factor (just to illustrate the use of extra parameters)
    metaprofiles(data.model) <- t(abs(pca$x[,1:factorization.rank])) / scale.factor

    # 3. Compute the offset as the mean expression
    data.model@offset <- rowMeans(x)

    # return updated data
    fit(seed) <- data.model
    seed
}

```

Then run the algorithm specifying it needs model `NMFOffset`:

```

# run custom algorithm with NMF model with offset
nmf(V, 3, my.algorithm.offset, model='NMFOffset', scale.factor=10)
<Object of class: NMFfit >
# Model:
  <Object of class: NMFOffset >
  features: 50
  basis/rank: 3
  samples: 20
  offset: [ 0.5363089 1.246813 0.5796618 1.591864 0.9041704 ... ]
# Details:
  algorithm: NMF.algo.f3f09d8
  seed: random
  distance metric: 'euclidean'
  residuals: 376.9509
  parameters:
  $scale.factor
  [1] 10

Timing:
   user  system elapsed
 0.000   0.000   0.005

```

3.2 Custom seeding method

The user can also define custom seeding method as a function of the form:

```

# start: object of class NMF
# target: the target matrix
my.seeding.method <- function(model, target){

  # use only the largest columns for W
  w.cols <- apply(target, 2, function(x) sqrt(sum(x^2)))
  metagenes(model) <- target[,order(w.cols)[1:nbasis(model)]]
}

```

```

    # initialize H randomly
    metaprofiles(model) <- matrix(runif(nbasis(model)*ncol(target))
                                   , nbasis(model), ncol(target))

    # return updated object
    return(model)
}

```

To use the new seeding method:

```

nmf(V, 3, 'snmf/r', seed=my.seeding.method)
<Object of class: NMFfit >
# Model:
<Object of class: NMFstd >
features: 50
basis/rank: 3
samples: 20
# Details:
algorithm: snmf/r
seed: NMF.seed.43d3bcd4
distance metric: 'euclidean'
residuals: 152.3523
Iterations: 75
Timing:
  user  system elapsed
 0.320   0.000   0.321

```

3.3 Package specific options

The package specific options can be retrieved or changed using the `nmf.getOption` and `nmf.options` functions. These behave similarly as the `getOption` and `options` base functions:

```

#show default algorithm and seeding method
nmf.options('default.algorithm', 'default.seed')
$default.algorithm
[1] "brunet"

$default.seed
[1] "random"

# retrieve a single option
nmf.getOption('default.seed')
[1] "random"

# All options
nmf.options()
$default.algorithm
[1] "brunet"

$track.interval
[1] 30

```

```

$error.track
[1] FALSE

$default.seed
[1] "random"

$parallel.backend
[1] "mc"

$verbose
[1] FALSE

$debug
[1] FALSE

```

References

- [Albright et al., 2006] R. Albright, J. Cox, D. Duling, A. Langville, C. Meyer (2006). Algorithms, initializations, and convergence for the nonnegative matrix factorization. *NCSU Technical Report Math 81706*. <http://meyer.math.ncsu.edu/Meyer/Abstracts/Publications.html>.
- [Badea L., 2008] Badea L. (2008). Profiles Common to Colon and Pancreatic Adenocarcinoma Using Simultaneous nonNegative. In *Pacific Symposium on Biocomputing, Volume 290* 2008:279–290.
- [Berry et al., 2006] Berry et al. (2006). Algorithms and Applications for Approximate Nonnegative Matrix Factorization. *Comput. Stat. Data Anal.*
- [Brunet et al., 2004] Brunet, J. P., Tamayo, P., Golub, T. R., and Mesirov, J. P. (2004). Metagenes and molecular pattern discovery using matrix factorization. *Proc Natl Acad Sci U S A*, **101**(12), 4164–4169.
- [A. Cichocki et al., 2004] Andrzej Cichocki , Rafal Zdunek, and Shun-ichi Amari (2004). New algorithms For Non-negative Matrix Factorization In Application To Blind Source Separation.
- [Chu et al., 2004] M.T. Chu, F. Diele, R. Plemmons, S. Ragni. Optimality, computation, and interpretation of nonnegative matrix factorizations. *Technical Report*, Departments of Mathematics and Computer Science, Wake Forest University, USA.
- [doMC, 2009] REvolution Computing: *doMC: Foreach parallel adaptor for the multicore package* 2009, [<http://CRAN.R-project.org/package=doMC>]. [R package version 1.2.0].
- [foreach, 2009] REvolution Computing: *doMC: Foreach parallel adaptor for the multicore package* 2009, [<http://CRAN.R-project.org/package=doMC>]. [R package version 1.3.0].
- [1] Frigyesi A, Höglund M: **Non-negative matrix factorization for the analysis of complex gene expression data: identification of clinically relevant tumor subtypes**. *Cancer informatics* 2008, **6**:275–292, [<http://view.ncbi.nlm.nih.gov/pubmed/19259414>].
- [2] Gao Y, Church G: **Improving molecular cancer class discovery through sparse non-negative matrix factorization**. *Bioinformatics* 2005, **21**(21):3970–3975, [<http://dx.doi.org/10.1093/bioinformatics/bti653>].

- [3] Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, Ellis B, Gautier L, Ge Y, Gentry J, Hornik K, Hothorn T, Huber W, Iacus S, Irizarry R, Leisch F, Li C, Maechler M, Rossini AJ, Sawitzki G, Smith C, Smyth G, Tierney L, Yang JYH, Zhang J: **Bioconductor: open software development for computational biology and bioinformatics.** *Genome biology* 2004, **5**:R80, [<http://www.ncbi.nlm.nih.gov/pubmed/15461798>].
- [4] Hutchins LNN, Murphy SMM, Singh P, Graber JHH: **Position-Dependent Motif Characterization Using Nonnegative Matrix Factorization.** *Bioinformatics (Oxford, England)* 2008, [<http://view.ncbi.nlm.nih.gov/pubmed/18852176>].
- [Kim and Park, 2007] Kim H, Park H: **Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis.** *Bioinformatics (Oxford, England)* 2007, **23**:1495–502, [<http://www.ncbi.nlm.nih.gov/pubmed/17483501>].
- [Lee and Seung, 2000] Lee, D.-D. and Seung, H.-S. (2000). Algorithms for non-negative matrix factorization. In *NIPS*, 556-562.
- [Pascual-Montano *et al.*, 2006] Pascual-Montano, A., Carazo, J.-M., Kochi, K., Lehmann, D., and Pascual-Marqui, R.-D. (2006). Nonsmooth nonnegative matrix factorization (nsnmf). *IEEE transactions on pattern analysis and machine intelligence*, **28**(3), 403-415.
- [R Software, 2008] R Development Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria. ISBN 3-900051-07-0. <http://www.R-project.org>.
- [5] Zhang J, Wei L, Feng X, Ma Z, Wang Y: **Pattern expression nonnegative matrix factorization: algorithm and applications to blind source separation.** *Computational intelligence and neuroscience* 2008, [<http://view.ncbi.nlm.nih.gov/pubmed/18566689>].