

# Generating heatmaps for Nonnegative Matrix Factorization

Package **NMF** - Version 0.6.1

Renaud Gaujoux

October 7, 2011

## Abstract

This vignette describes how to produce different informative heatmaps from NMF objects, such as returned by the function **nmf** in the package **NMF** [?]. The main drawing engine is implemented by the function **aheatmap**, which is a modification of the function **pheatmap** from the package **pheatmap** [?], and provides convenient and quick ways of producing high quality and customizable annotated heatmaps. Currently this function is part of the package **NMF**, but may eventually compose a separate package by itself.

## Contents

<b>1 Preliminaries</b>	<b>1</b>	<b>3 Basis matrix: basismap</b>	<b>5</b>
1.1 Quick reminder on NMF models . . .	1	<b>4 Consensus matrix: consensusmap</b>	<b>6</b>
1.2 Heatmaps for NMF . . . . .	1	4.1 Single fit . . . . .	7
1.3 Heatmap engine . . . . .	2	4.2 Single method over a range of ranks	8
1.4 Data and model . . . . .	2	4.3 Single rank over a range of methods	9
<b>2 Mixture Coefficient matrix: coefmap</b>	<b>3</b>	<b>5 Generic heatmap engine: aheatmap</b>	<b>9</b>

## 1 Preliminaries

### 1.1 Quick reminder on NMF models

Given a nonnegative target matrix  $X$  of dimension  $n \times p$ , NMF algorithms aim at finding a rank  $k$  approximation of the form:

$$X \approx WH,$$

where  $W$  and  $H$  are nonnegative matrices of dimensions  $n \times k$  and  $k \times p$  respectively.

The matrix  $W$  is the basis matrix, whose columns are the basis components. The matrix  $H$  is the mixture coefficient or weight matrix, whose columns contain the contribution of each basis component to the corresponding column of  $X$ . We call the rows of  $H$  the basis profiles.

### 1.2 Heatmaps for NMF

Because NMF objects essentially wrap up a pair of matrices, heatmaps are convenient to visualise the results of NMF runs. The package **NMF** provides several specialised heatmap functions, designed to produce heatmaps with sensible default configurations according to the data being drawn. Being all based on a common drawing engine they share almost identical interfaces and capabilities. The following specialised functions are currently implemented:

**basismap** draws heatmaps of the basis matrix

**coefmap** draws heatmaps of the mixture coefficient matrix

**consensusmap** draws heatmaps of the consensus matrix, for results of multiple NMF runs.

### 1.3 Heatmap engine

All the above functions eventually call a common heatmap engine, with different default parameters, chosen to be relevant for the given underlying data. The engine is implemented by the function `aheatmap`. Its development started as modification of the function `pheatmap` from the package `pheatmap` [?]. The initial objective was to improve and increase its capabilities, as well as a defining a simplified interface, more consistent with the R core function `heatmap`. We eventually aim at providing a general, flexible, powerful and easy to use engine for drawing annotated heatmaps.

The function `aheatmap` has many advantages compared to other heatmap functions such as `heatmap`, `gplots::heatmap2`, `heatmap.plus::heatmap.plus`, or even `pheatmap`:

- Annotations: unlimited number of annotation tracks can be added to *both* columns and rows, with automated colouring for categorical and numeric variables.
- Compatibility with both base and grid graphics: the function can be directly called in drawing contexts such as `grid`, `mfrow` or `layout`. This is a feature many R users were looking for, and that was strictly impossible with base heatmaps.
- Legends: default automatic legend and colouring;
- Customisation: clustering methods, annotations, colours and legend can all be customised, even separately for rows and columns;
- Convenient interface: many arguments provide multiple ways of specifying their value(s), which speeds up developping/writing and reduce the amount of code required to generate customised plots (e.g. see ??).
- Aesthetics: the heatmaps look globally cleaner, the image and text components are by default well proportioned relatively to each other, and all fit within the graphic device.

### 1.4 Data and model

For the purpose of illustrating the use of each heatmap function, we generate a random target matrix, as well as some annotations or covariates:

```
# random data matrix
X <- rmatrix(100,20)
# row annotations and covariates
n <- nrow(X)
d <- rnorm(n)
e <- unlist(mapply(rep, c('Brain', 'Muscle', 'Skin'), 10))
e <- c(e, rep(NA, n-length(e)))
# column annotations and covariates
p <- ncol(X)
a <- sample(c('alpha', 'beta', 'gamma'), p, replace=TRUE)
b <- sample(letters[10:15], p, replace=TRUE)
c <- rnorm(p)
# gather them in a data.frame
covariates <- data.frame(a, b, c)
```

Then, we fit an NMF model using multiple runs:

```
res <- nmf(X, 3, nrun=10)

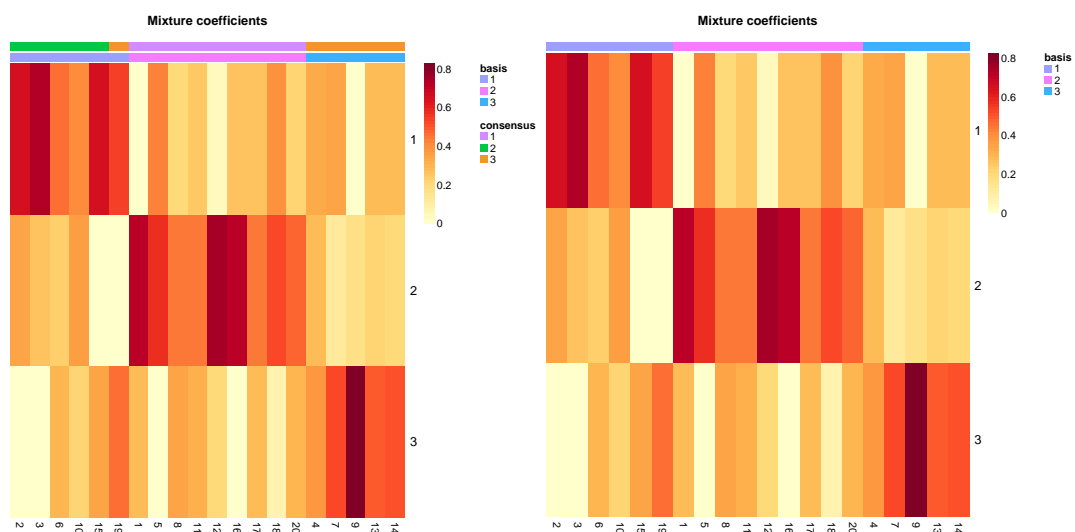
<Object of class: NMFfitX1 >
  Method: brunet
  Runs: 10
  RNG:
    rstream.mrg32k3a - 931979, 773177, 849931, 312446, 81019, 29408
    [0ae4de264fbe988e14294c8986ab28e7]
  Total timing:
    user system elapsed
    3.610 0.060 7.019
```

**NB:** To keep the vignette simple, we always use the default NMF method (i.e. 'brunet'), but all steps could be performed using a different method, or multiple methods in order to compare their performances.

## 2 Mixture Coefficient matrix: coefmap

The basis matrix of the result can be plotted using the function `coefmap`. The default behaviour for multiple NMF runs is to add two annotation tracks that show the clusters obtained by the best fit and the hierarchical clustering of the consensus matrix<sup>1</sup>. In the legend, these tracks are named *basis* and *consensus* respectively. For single NMF run or NMF model objects, no consensus data are available, and only the clusters from the fit are displayed.

```
opar <- par(mfrow=c(1,2))
# coefmap from multiple run fit: includes a consensus track
coefmap(res)
# coefmap of a single run fit: no consensus track
coefmap(minfit(res))
par(opar)
```



<sup>1</sup>The hierarchical clustering is computed using the consensus matrix itself as a similarity measure, and average linkage. See `?consensusmap`.

**NB:** Note how both heatmaps were drawn on the same plot, simply using the standard call to `par(mfrow=c(1,2))`. This is impossible to achieve with the R core function `heatmap`. See section 5 for more details about compatibility with base and grid graphics.

By default:

- the rows are not ordered;
- the columns are grouped by the clusters defined by the dominant basis component for each column;
- each column is scaled to sum up to one;
- the color palette used is 'YlOrRd' from the package `RColorBrewer` [?].

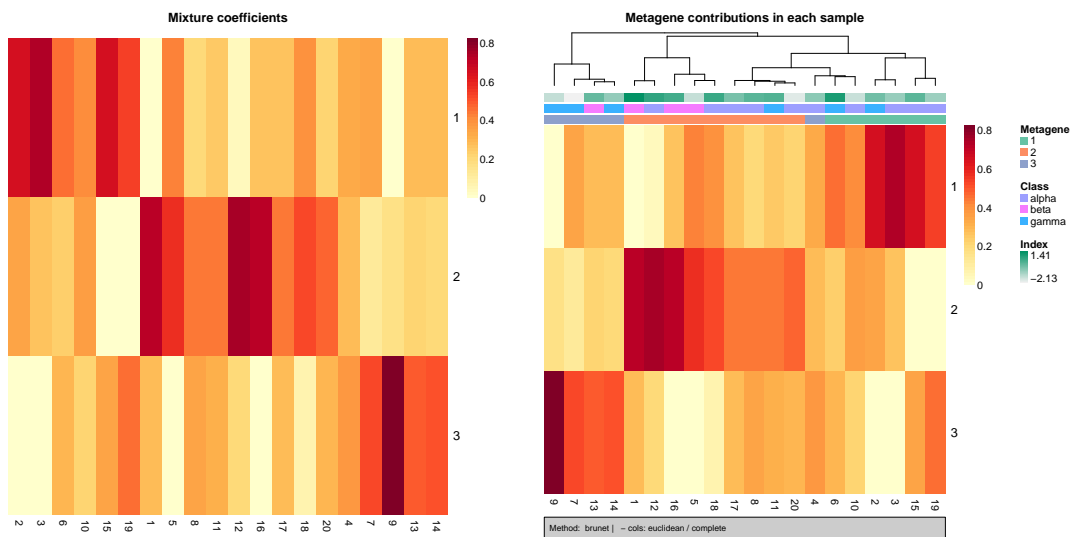
In term of arguments passed to the heatmap engine `aheatmap`, these default settings translate as:

```
Rowv = NA
Colv = order(as.numeric(predict(object)))
scale = 'c1'
color = 'YlOrRd'
```

Since the ordering does not come from hierarchical clustering, no dendrogram is displayed. The default behaviour of `aheatmap` can be obtained by setting arguments `Rowv=TRUE`, `Colv=TRUE`, `scale='none'`.

The automatic annotation tracks can be hidden all together by setting argument `tracks=NA`, displayed separately by passing only one of the given names (e.g. `tracks='basis'`), and their legend names may be changed by specifying e.g. `tracks=c(Metagene='basis', 'consensus')`. Beside this, they are handled by the heatmap engine function `aheatmap` and can be customised as any other annotation tracks – that can be added via argument `annCol` (see 5 or `?aheatmap` for more details).

```
opar <- par(mfrow=c(1,2))
# removing all automatic annotation tracks
coefmap(res, tracks=NA)
# customized plot
coefmap(res, Colv = 'euclidean'
        , main = "Metagene contributions in each sample", labCol = NULL
        , tracks = c(Metagene='basis'), annCol = list(Class=a, Index=c)
        , annColors = list(Metagene='Set2')
        , info = TRUE)
par(opar)
```



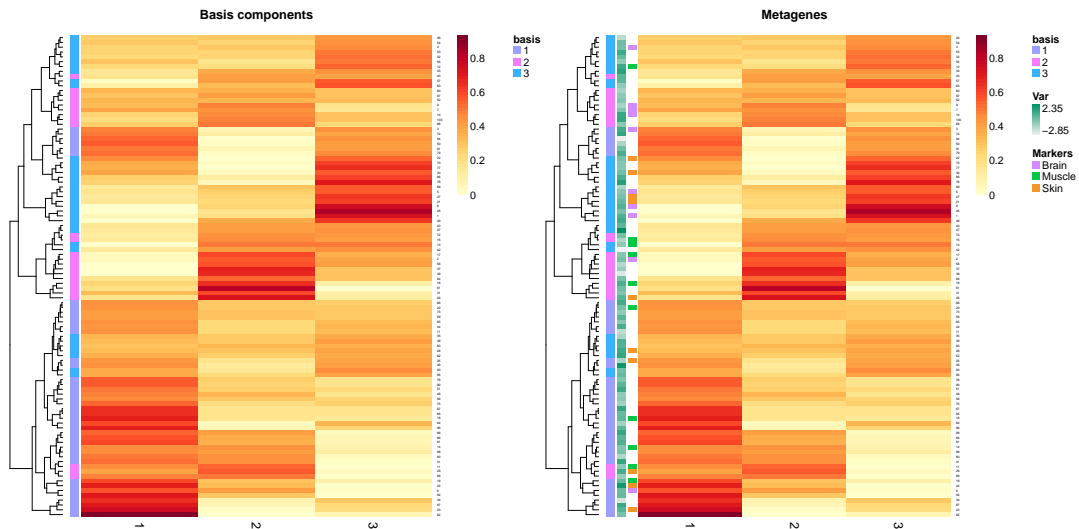
**NB:** The feature that allows to display some information about the fit at the bottom of the plot via argument `info=TRUE` is still experimental. It is helpful mostly when developing algorithms or doing an analysis, but would seldom be used in publications.

### 3 Basis matrix: `basismap`

The basis matrix can be plotted using the function `basismap`. The default behaviour is to add an annotation track that shows for each row the dominant basis component. That is, for each row, the index of the basis component with the highest loading.

Again, the track can be disabled by setting `tracks=NA`, and extra row annotations can be added using argument `annRow`.

```
opar <- par(mfrow=c(1,2))
# default plot
basismap(res)
# customized plot
basismap(res, main="Metagenes", annRow=list(Var=d, Markers=e))
par(opar)
```



By default:

- the columns are not ordered;
- the rows are ordered by hierarchical clustering using default distance and linkage methods ('euclidean' and 'complete');
- each row is scaled to sum up to one;
- the color palette used is 'YlOrRd' from the package RColorBrewer [?].

In term of arguments passed to the heatmap engine `aheatmap`, these default settings translate as:

```
Colv = NA
scale = 'r1'
color = 'YlOrRd'
```

## 4 Consensus matrix: consensusmap

When doing clustering with NMF, a common way of assessing the stability of the clusters obtained for a given rank is to consider the consensus matrix computed over multiple independent NMF runs, which is the average of the connectivity matrices of each separate run <sup>2</sup>. This procedure is usually repeated over a certain range of factorization ranks, and the results are compared to identify which rank gives the best clusters, possibly in the light of some extra knowledge one could have about the samples (e.g. covariates). The functions `nmf` and `consensusmap` make it easy to implement this whole process.

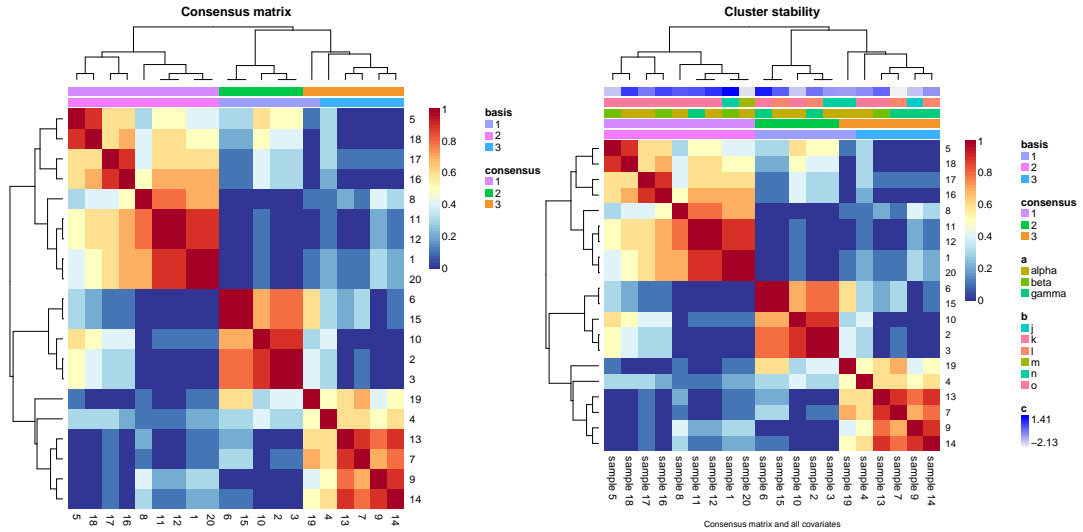
**NB:** The consensus plots can also be generated for fits obtained from single NMF runs, in which case the consensus matrix simply reduces to a single connectivity matrix. This is a binary matrix (i.e. entries are either 0 or 1), that will always produce a bi-colour heatmap, and by default clear blocks for each cluster.

<sup>2</sup>Hence, stability here means robustness with regards to the initial starting point, and shall not be interpreted as in e.g. cross-validation/bootstrap analysis. However, one can argue that having very consistent clusters across runs somehow supports for a certain regularity or the presence of an underlying pattern in the data.

## 4.1 Single fit

In section 1.4, the NMF fit `res` was computed with argument `nrun=10`, and therefore contains the best fit over 10 runs, as well as the consensus matrix computed over all the runs<sup>3</sup>. This can be plotted using the function `consensusmap`, which allows for the same kind of customization as the other NMF heatmap functions:

```
opar <- par(mfrow=c(1,2))
# default plot
consensusmap(res)
# customized plot
consensusmap(res, annCol=covariates, annColors=list(c='blue')
              , labCol='sample ', main='Cluster stability'
              , sub='Consensus matrix and all covariates')
par(opar)
```



By default:

- the rows and columns of the consensus heatmap are symmetrically ordered by hierarchical clustering using the consensus matrix as a similarity measure and average linkage, and the associated dendrogram is displayed;
- the color palette used is the reverse of 'RdYlBu' from the package RColorBrewer [?].

In term of arguments passed to the heatmap engine `aheatmap`, these default settings translate as:

```
distfun = function(x) as.dist(1-x) # x being the consensus matrix
hclustfun = 'average'
Rowv = TRUE
Colv = "Rowv"
color = '-RdYlBu'
```

<sup>3</sup>If one were interested in keeping the fits from all the runs, the function `nmf` should have been called with argument `.options='k'`. See section *Options* in `?nmf`. The downstream handling of the result would remain identical.

## 4.2 Single method over a range of ranks

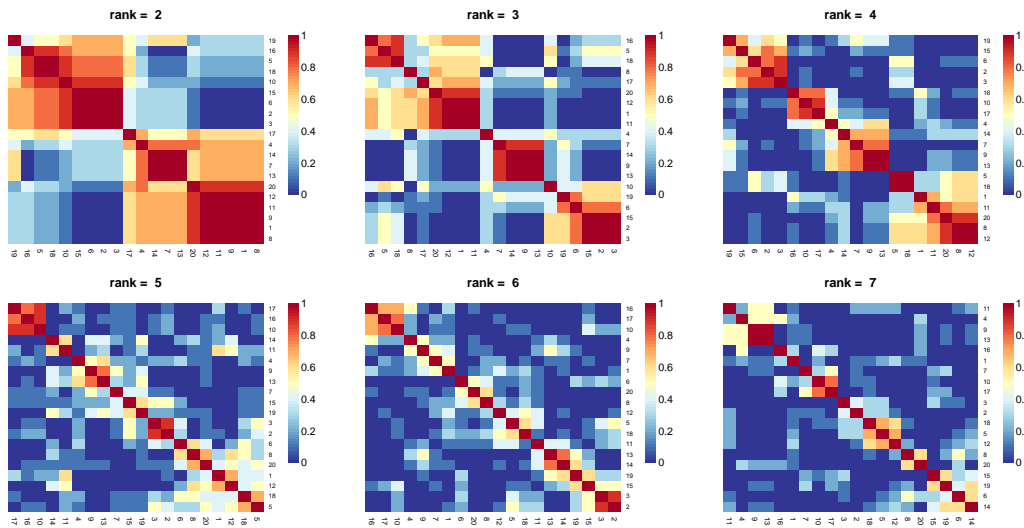
The function `nmf` accepts a range of value for the rank (argument `rank`), making it fit NMF models for each value in the given range <sup>4</sup>:

```
res2_7 <- nmf(X, 2:7, nrun=10, .options='v')
Compute NMF rank= 2 ... + measures ... OK
Compute NMF rank= 3 ... + measures ... OK
Compute NMF rank= 4 ... + measures ... OK
Compute NMF rank= 5 ... + measures ... OK
Compute NMF rank= 6 ... + measures ... OK
Compute NMF rank= 7 ... + measures ... OK

class(res2_7)
[1] "NMF.rank"
```

The result `res2_7` is an S3 object of class `'NMF.rank'`, that contains – amongst other data – a list of the best fits obtained for each value of the rank in range `[[2,7]]`. The method of `consensusmap` defined for class `'NMF.rank'`, which plots all the consensus matrices on the same plot:

```
consensusmap(res2_7)
```



**NB:** The main title of each consensus heatmap can be customized by passing to argument `main` a character vector or a list whose elements specify each title. All other arguments are used in each internal call to `consensusmap`, and will therefore affect all the plots simultaneously. The layout can be specified via argument `layout` as a numeric vector giving the number of rows and columns in a `mrow`-like way, or as a matrix that will be passed to R core function `layout`. See `?consensusmap` for more details and example code.

<sup>4</sup>Before version 0.6, this feature was provided by the function `nmfEstimateRank`. From version 0.6, the function `nmf` accepts ranges of ranks, and internally calls the function `nmfEstimateRank` – that remains exported and can still be called directly. See documentation `?nmfEstimateRank` for more details on the returned value.



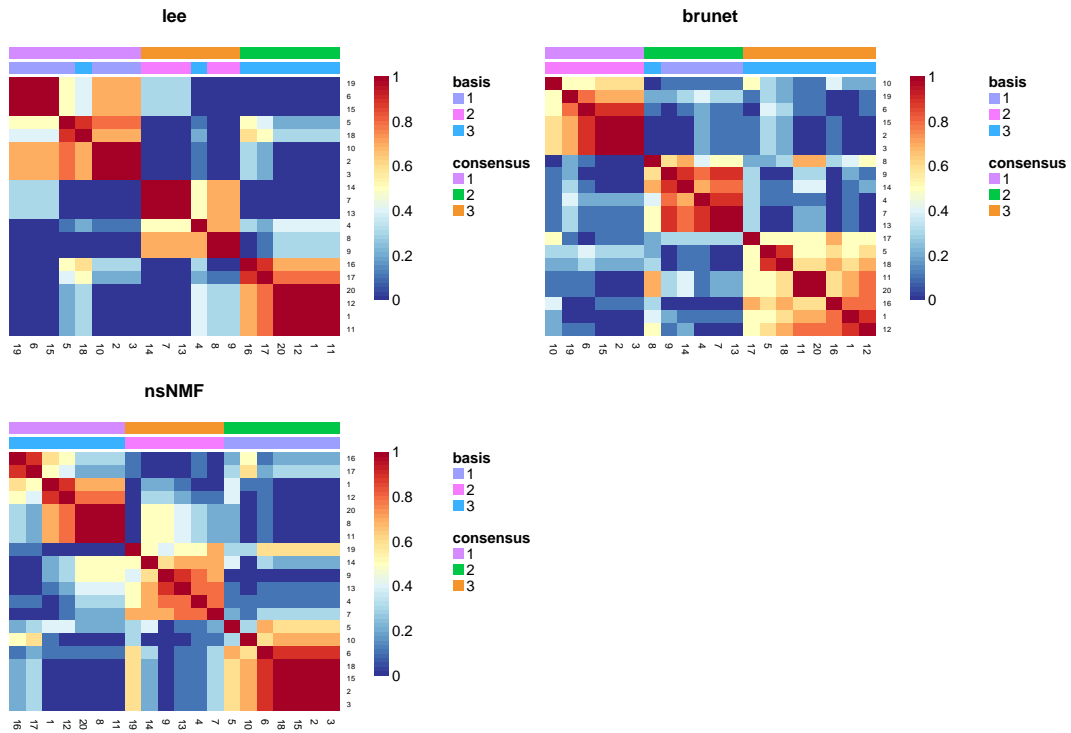
### 4.3 Single rank over a range of methods

If one is interested in comparing methods, for a given factorization rank, then one can fit an NMF model for each method by providing the function `nmf` with a list in argument `method`:

```
res_methods <- nmf(X, 3, list('lee', 'brunet', 'nsNMF'), nrun=10)
class(res_methods)
[1] "NMFList"
attr("package")
[1] "NMF"
```

The result `res_methods` is an S4 object of class `NMFList`, which is essentially a named list, that contains each fits and the CPU time required by the whole computation. As previously, the sequence of consensus matrices is plotted with `consensusmap`:

```
consensusmap(res_methods)
```



## 5 Generic heatmap engine: aheatmap