# **Software Manual**





# **APCluster**

# An R Package for Affinity Propagation Clustering

Ulrich Bodenhofer, Johannes Palme, Chrats Melkonian, and Andreas Kothmeier

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

Version 1.4.3, February 24, 2016

# **Scope and Purpose of this Document**

This document is a user manual for the R package apcluster [1]. It is only meant as a gentle introduction into how to use the basic functions implemented in this package. Not all features of the R package are described in full detail. Such details can be obtained from the documentation enclosed in the R package. Further note the following: (1) this is neither an introduction to affinity propagation nor to clustering in general; (2) this is not an introduction to R. If you lack the background for understanding this manual, you first have to read introductory literature on these subjects.

Contents 3

$\boldsymbol{\cap}$	4 .	4
Con	ITET	) TC

1	Introduction	4
2	Installation2.1 Installation via CRAN2.2 Manual installation from source2.3 Compatibility issues	4 4 4 5
3	Getting Started	5
4	Adjusting Input Preferences	
5	Exemplar-based Agglomerative Clustering  5.1 Getting started	20 20 23 26
6	Leveraged Affinity Propagation	27
7	Sparse Affinity Propagation	30
8	B A Toy Example with Biological Sequences	
9	Similarity Matrices  9.1 The function negDistMat()	39 39 42 45 46 48
10	Miscellaneous  10.1 Convenience vs. efficiency  10.2 Clustering named objects  10.3 Computing a label vector from a clustering result  10.4 Customizing heatmaps  10.5 Adding a legend to plots of clustering results  10.6 Implementation and performance issues  10.7 Using APCluster in Conjunction with KeBABS	49 50 51 52 54 55 56
11	Special Notes for Users Upgrading from Previous Versions  11.1 Upgrading from a version older than 1.3.0	<b>56</b> 56 56
12	Change Log	57
13	How to Cite This Package	60

4 2 Installation

#### 1 Introduction

Affinity propagation (AP) is a relatively new clustering algorithm that has been introduced by Brendan J. Frey and Delbert Dueck [5]. The authors themselves describe affinity propagation as follows:<sup>2</sup>

"An algorithm that identifies exemplars among data points and forms clusters of data points around these exemplars. It operates by simultaneously considering all data point as potential exemplars and exchanging messages between data points until a good set of exemplars and clusters emerges."

AP has been applied in various fields recently, among which bioinformatics is becoming increasingly important. Frey and Dueck have made their algorithm available as Matlab code. Matlab, however, is relatively uncommon in bioinformatics. Instead, the statistical computing platform R has become a widely accepted standard in this field. In order to leverage affinity propagation for bioinformatics applications, we have implemented affinity propagation as an R package. Note, however, that the given package is in no way restricted to bioinformatics applications. It is as generally applicable as Frey's and Dueck's original Matlab code. I

Starting with Version 1.1.0, the apcluster package also features exemplar-based agglomerative clustering which can be used as a clustering method on its own or for creating a hierarchy of clusters that have been computed previously by affinity propagation. Leveraged Affinity Propagation, a variant of AP especially geared to applications involving large data sets, has first been included in Version 1.3.0.

#### 2 Installation

#### 2.1 Installation via CRAN

The R package apcluster (current version: 1.4.3) is part of the *Comprehensive R Archive Network (CRAN)*<sup>3</sup>. The simplest way to install the package, therefore, is to enter the following command into your R session:

```
install.packages("apcluster")
```

If you use R on Windows or Mac OS, you can also conveniently use the package installation menu of your R GUI.

#### 2.2 Manual installation from source

Under special circumstances, e.g. if you want to compile the C++ code included in the package with some custom options, you may prefer to install the package manually from source. To this end, open the package's page at CRAN<sup>4</sup> and then proceed as follows:

<sup>1</sup>http://www.psi.toronto.edu/affinitypropagation/

 $<sup>^2</sup> quoted\ from\ http://www.psi.toronto.edu/affinitypropagation/faq.html\#def$ 

<sup>3</sup>http://cran.r-project.org/

<sup>4</sup>http://cran.r-project.org/web/packages/apcluster/index.html

3 Getting Started

- 1. Download apcluster\_1.4.3.tar.gz and save it to your harddisk.
- 2. Open a shell/terminal/command prompt window and change to the directory where you put apcluster\_1.4.3.tar.gz. Enter

```
R CMD INSTALL apcluster_1.4.3.tar.gz
```

to install the package.

Note that this might require additional software on some platforms. Windows requires Rtools<sup>5</sup> to be installed and to be available in the default search path (environment variable PATH). Mac OS X requires Xcode developer tools<sup>6</sup> (make sure that you have the command line tools installed with Xcode).

#### 2.3 Compatibility issues

All versions downloadable from CRAN have been built using the latest version, R 3.2.3. However, the package should work without severe problems on R versions  $\geq$ 3.0.0.

#### 3 Getting Started

To load the package, enter the following in your R session:

```
library(apcluster)
```

If this command terminates without any error message or warning, you can be sure that the package has been installed successfully. If so, the package is ready for use now and you can start clustering your data with affinity propagation.

The package includes both a user manual (this document) and a reference manual (help pages for each function). To view the user manual, enter

```
vignette("apcluster")
```

Help pages can be viewed using the help command. It is recommended to start with

```
help(apcluster)
```

Affinity propagation does not require the data samples to be of any specific kind or structure. AP only requires a *similarity matrix*, i.e., given l data samples, this is an  $l \times l$  real-valued matrix S, in which an entry  $S_{ij}$  corresponds to a value measuring how similar sample i is to sample j. AP does not require these values to be in a specific range. Values can be positive or negative. AP does not even require the similarity matrix to be symmetric (although, in most applications, it will

<sup>5</sup>http://cran.r-project.org/bin/windows/Rtools/

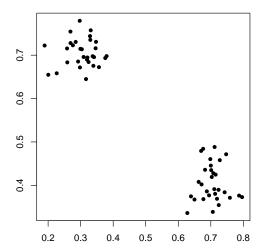
<sup>6</sup>https://developer.apple.com/technologies/tools/

6 3 Getting Started

be symmetric anyway). A value of  $-\infty$  is interpreted as "absolute dissimilarity". The higher a value, the more similar two samples are considered.

To get a first impression, let us create a random data set in  $\mathbb{R}^2$  as the union of two "Gaussian clouds":

```
cl1 <- cbind(rnorm(30, 0.3, 0.05), rnorm(30, 0.7, 0.04))
cl2 <- cbind(rnorm(30, 0.7, 0.04), rnorm(30, 0.4, .05))
x1 <- rbind(cl1, cl2)
plot(x1, xlab="", ylab="", pch=19, cex=0.8)</pre>
```



The package apcluster offers several different ways for clustering data. The simplest way is the following:

```
apres1a <- apcluster(negDistMat(r=2), x1)</pre>
```

In this example, the function apcluster() first computes a similarity matrix for the input data x1 using the *similarity function* passed as first argument. The choice negDistMat(r=2) is the standard similarity measure used in the papers of Frey and Dueck — negative squared distances.

Alternatively, one can compute the similarity matrix beforehand and call apcluster() for the similarity matrix (for a more detailed description of the differences, see 10.1):

```
s1 <- negDistMat(x1, r=2)
apres1b <- apcluster(s1)</pre>
```

The function apcluster() creates an object belonging to the S4 class APResult which is defined by the present package. To get detailed information on which data are stored in such objects, enter

3 Getting Started

```
help(APResult)
```

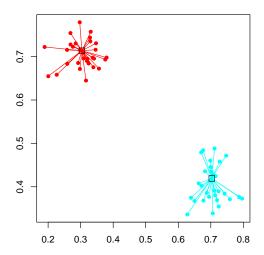
The simplest thing we can do is to enter the name of the object (which implicitly calls show()) to get a summary of the clustering result:

```
apres1a
##
## APResult object
##
## Number of samples = 60
## Number of iterations = 130
## Input preference = -0.1582302
## Sum of similarities = -0.1974598
## Sum of preferences = -0.3164603
## Net similarity = -0.5139202
## Number of clusters = 2
##
## Exemplars:
    23 49
##
## Clusters:
   Cluster 1, exemplar 23:
##
       1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
##
        26 27 28 29 30
##
     Cluster 2, exemplar 49:
##
        31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52
##
      53 54 55 56 57 58 59 60
```

The appluster package allows for plotting the original data set along with a clustering result:

```
plot(apres1a, x1)
```

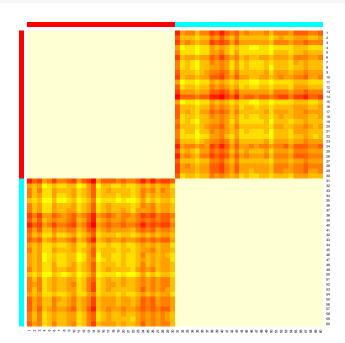
8 3 Getting Started



In this plot, each color corresponds to one cluster. The exemplar of each cluster is marked by a box and all cluster members are connected to their exemplars with lines.

A heatmap is plotted with heatmap():

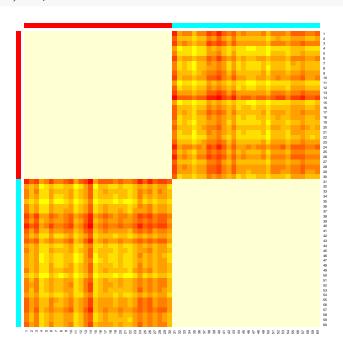
#### heatmap(apres1a)



In the heatmap, the samples are grouped according to clusters. The above heatmap confirms again that there are two main clusters in the data. A heatmap can be plotted for the object apres1a be-

cause apcluster(), if called for data and a similarity function, by default includes the similarity matrix in the output object (unless it was called with the switch includeSim=FALSE). If the similarity matrix is not included (which is the default if apcluster() has been called on a similarity matrix directly), heatmap() must be called with the similarity matrix as second argument:

heatmap(apres1b, s1)

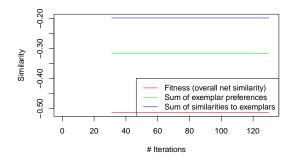


Suppose we want to have better insight into what the algorithm did in each iteration. For this purpose, we can supply the option details=TRUE to apcluster():

```
apres1c <- apcluster(s1, details=TRUE)</pre>
```

This option tells the algorithm to keep a detailed log about its progress. For example, this allows for plotting the three performance measures that AP uses internally for each iteration:

```
plot(apres1c)
```



These performance measures are:

- 1. Sum of exemplar preferences
- 2. Sum of similarities of exemplars to their cluster members
- 3. Net fitness: sum of the two former

For details, the user is referred to the original affinity propagation paper [5] and the supplementary material published on the affinity propagation Web page. We see from the above plot that the algorithm has not made any change for the last 100 (of 130!) iterations. AP, through its parameter convits, allows to control for how long AP waits for a change until it terminates (the default is convits=100). If the user has the feeling that AP will probably converge quicker on his/her data set, a lower value can be used:

```
apres1c <- apcluster(s1, convits=15, details=TRUE)
apres1c
##
## APResult object
##
## Number of samples
                            60
## Number of iterations = 45
## Input preference
                       = -0.1582302
## Sum of similarities = -0.1974598
## Sum of preferences
                         =
                           -0.3164603
## Net similarity
                         = -0.5139202
## Number of clusters
##
## Exemplars:
     23 49
##
## Clusters:
##
     Cluster 1, exemplar 23:
        1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
##
         26 27 28 29 30
     Cluster 2, exemplar 49:
##
        31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52
##
         53 54 55 56 57 58 59 60
```

# 4 Adjusting Input Preferences

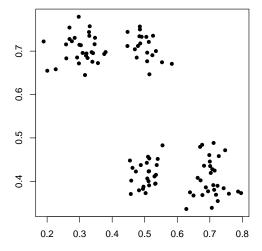
Apart from the similarity matrix itself, the most important input parameter of AP is the so-called *input preference* which can be interpreted as the tendency of a data sample to become an exemplar (see [5] and supplementary material on the AP homepage<sup>1</sup> for a more detailed explanation). This input preference can either be chosen individually for each data sample or it can be a single value

shared among all data samples. Input preferences largely determine the number of clusters, in other words, how fine- or coarse-grained the clustering result will be.

The input preferences one can specify for AP are roughly in the same range as the similarity values, but they do not have a straightforward interpretation. Frey and Dueck have introduced the following rule of thumb: "The shared value could be the median of the input similarities (resulting in a moderate number of clusters) or their minimum (resulting in a small number of clusters)." [5]

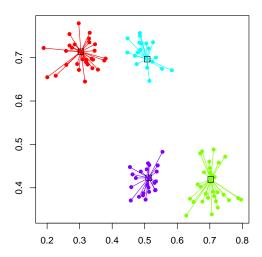
Our AP implementation uses the median rule by default if the user does not supply a custom value for the input preferences. In order to provide the user with a knob that is — at least to some extent — interpretable, the function apcluster() provides an argument q that allows to set the input preference to a certain quantile of the input similarities: resulting in the median for q=0.5 and in the minimum for q=0. As an example, let us add two more "clouds" to the data set from above:

```
cl3 <- cbind(rnorm(20, 0.50, 0.03), rnorm(20, 0.72, 0.03))
cl4 <- cbind(rnorm(25, 0.50, 0.03), rnorm(25, 0.42, 0.04))
x2 <- rbind(x1, cl3, cl4)
plot(x2, xlab="", ylab="", pch=19, cex=0.8)</pre>
```



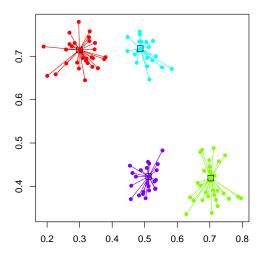
For the default setting, we obtain the following result:

```
apres2a <- apcluster(negDistMat(r=2), x2)
plot(apres2a, x2)</pre>
```



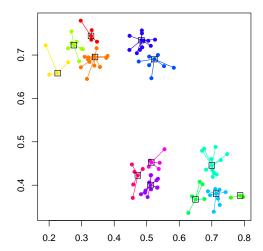
For the minimum of input similarities, we obtain the following result:

```
apres2b <- apcluster(negDistMat(r=2), x2, q=0)
plot(apres2b, x2)</pre>
```



So we see that AP is quite robust against a reduction of input preferences in this example which may be caused by the clear separation of the four clusters. If we increase input preferences, however, we can force AP to split the four clusters into smaller sub-clusters:

```
apres2c <- apcluster(negDistMat(r=2), x2, q=0.8)
plot(apres2c, x2)</pre>
```

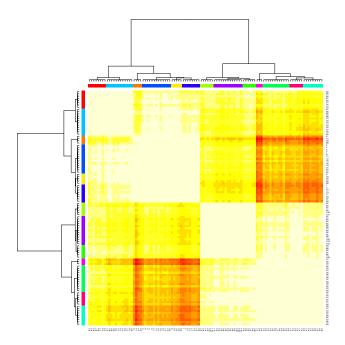


Note that the input preference used by AP can be recovered from the output object (no matter which method to adjust input preferences has been used). On the one hand, the value is printed if the object is displayed (by show or by entering the output object's name). On the other hand, the value can be accessed directly via the slot p:

```
apres2c@p
## [1] -0.008658507
```

As noted above already, we can produce a heatmap by calling heatmap() for an APResult object:

```
heatmap(apres2c)
```



The order in which the clusters are arranged in the heatmap is determined by means of joining the cluster agglomeratively (see Section 5 below). Although the affinity propagation result contains 13 clusters, the heatmap indicates that there are actually four clusters which can be seen as very brightly colored squares along the diagonal. We also see that there seem to be two pairs of adjacent clusters, which can be seen from the fact that there are two relatively light-colored blocks along the diagonal encompassing two of the four clusters in each case. If we look back at how the data have been created (see also plots above), this is exactly what is to be expected.

The above example with q=0 demonstrates that setting input preferences to the minimum of input similarities does not necessarily result in a very small number of clusters (like one or two). This is due to the fact that input preferences need not necessarily be exactly in the range of the similarities. To determine a meaningful range, an auxiliary function is available which, in line with Frey's and Dueck's Matlab code, allows to compute a minimum value (for which one or at most two clusters would be obtained) and a maximum value (for which as many clusters as data samples would be obtained):

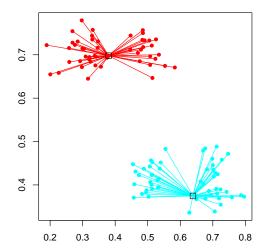
```
preferenceRange(apres2b@sim)
## [1] -5.717876e+00 -1.596955e-06
```

The function returns a two-element vector with the minimum value as first and the maximum value as second entry. The computations are done approximately by default. If one is interested in exact bounds, supply exact=TRUE (resulting in longer computation times).

Many clustering algorithms need to know a pre-defined number of clusters. This is often a major nuisance, since the exact number of clusters is hard to know for non-trivial (in particular, high-dimensional) data sets. AP avoids this problem. If, however, one still wants to require a fixed

number of clusters, this has to be accomplished by a search algorithm that adjusts input preferences in order to produce the desired number of clusters in the end. For convenience, this search algorithm is available as a function apclusterK() (analogous to Frey's and Dueck's Matlab implementation<sup>1</sup>). We can use this function to force AP to produce only two clusters (merging the two pairs of adjacent clouds into one cluster each). Analogously to apcluster(), apclusterK() supports two variants — it can either be called for a similarity measure and data or on a similarity matrix directly.

```
apres2d <- apclusterK(negDistMat(r=2), x2, K=2, verbose=TRUE)
## Trying p = -0.005719472
##
      Number of clusters: 14
## Trying p = -0.05718034
##
      Number of clusters: 4
## Trying p = -0.5717891
##
      Number of clusters: 3
## Trying p = -2.858939 (bisection step no. 1)
##
      Number of clusters: 2
##
## Number of clusters: 2 \text{ for p} = -2.858939
plot(apres2d, x2)
```

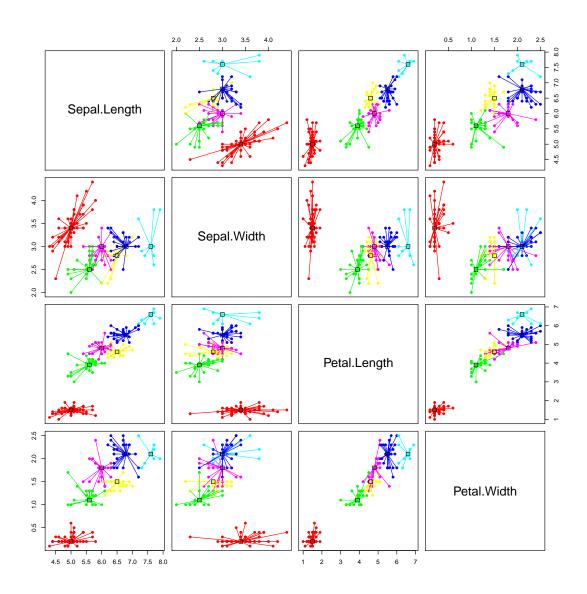


Now let us quickly consider a simple data set with more than two features. The notorious example is Fisher's iris data set:

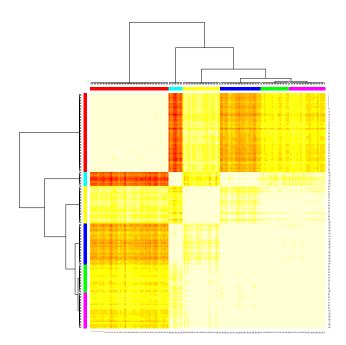
```
data(iris)
apIris1 <- apcluster(negDistMat(r=2), iris)</pre>
apIris1
##
## APResult object
##
## Number of samples
                       = 150
## Number of iterations = 162
## Input preference = -5.57
## Sum of similarities = -45.96
## Sum of preferences = -33.42
## Net similarity
                    = -79.38
## Number of clusters
##
## Exemplars:
     8 55 70 106 113 139
##
## Clusters:
##
     Cluster 1, exemplar 8:
##
        1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
        26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
##
        48 49 50
##
     Cluster 2, exemplar 55:
        51 52 53 55 59 66 69 73 74 75 76 77 78 87 88 98 120 134
##
##
     Cluster 3, exemplar 70:
##
        54 58 60 61 63 65 68 70 72 80 81 82 83 89 90 91 93 94 95 96 97 99
##
        100 107
     Cluster 4, exemplar 106:
##
##
        106 108 110 118 119 123 131 132 136
##
     Cluster 5, exemplar 113:
##
        101 103 104 105 109 111 112 113 116 117 121 125 126 129 130 133
         137 138 140 141 142 144 145 146 148 149
##
##
     Cluster 6, exemplar 139:
        56 57 62 64 67 71 79 84 85 86 92 102 114 115 122 124 127 128 135
##
##
        139 143 147 150
```

AP has identified 6 clusters. Since Version 1.3.2, the package also allows for superimposing clustering results in scatter plot matrices:

```
plot(apIris1, iris)
```



The heatmap looks as follows:



Now let us try to obtain fewer clusters by using the minimum of off-diagonal similarities:

```
data(iris)
apIris2 <- apcluster(negDistMat(r=2), iris, q=0)</pre>
apIris2
##
## APResult object
##
## Number of samples = 150
## Number of iterations = 126
## Input preference = -50.2
## Sum of similarities = -84.44
## Sum of preferences = -150.6
## Net similarity = -235.04
## Number of clusters
##
## Exemplars:
##
     8 56 113
## Clusters:
     Cluster 1, exemplar 8:
##
##
        1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
        26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
        48 49 50
##
##
     Cluster 2, exemplar 56:
  52 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74
##
```

```
## 75 76 79 80 81 82 83 84 85 86 88 89 90 91 92 93 94 95 96 97 98 99

## 100 102 107 114 120 122 124 127 128 134 139 143 150

## Cluster 3, exemplar 113:

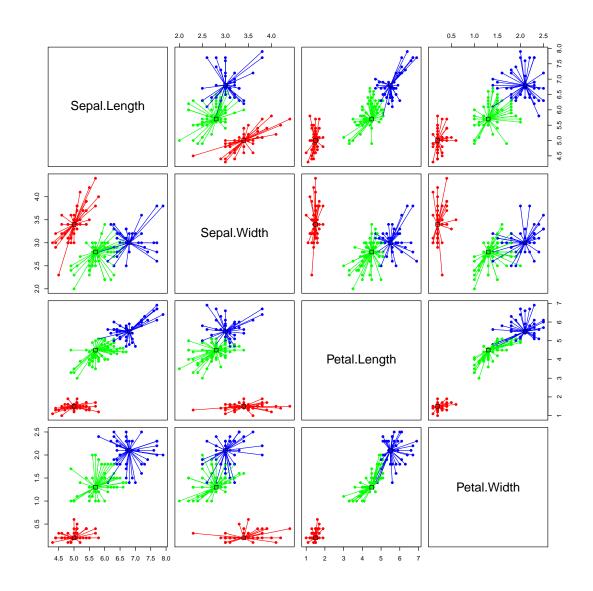
## 51 53 77 78 87 101 103 104 105 106 108 109 110 111 112 113 115

## 116 117 118 119 121 123 125 126 129 130 131 132 133 135 136 137

## 138 140 141 142 144 145 146 147 148 149
```

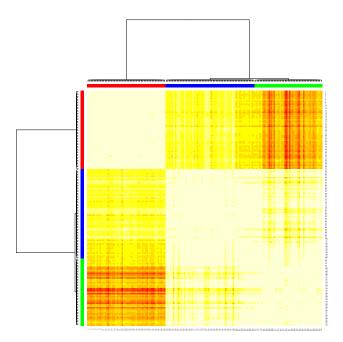
AP has identified 3 clusters. If we again superimpose them in the scatter plot matrix, we obtain the following:

```
plot(apIris2, iris)
```



Finally, the heatmap looks as follows:





So, looking at the heatmap, the 3 clusters seem quite reasonable, at least in the light of the fact that there are three species in the data set, *Iris setosa*, *Iris versicolor*, and *Iris virginica*, where *Iris setosa* is very clearly separated from each other (first cluster in the heatmap) and the two others are partly overlapping.

## 5 Exemplar-based Agglomerative Clustering

The function aggExCluster() realizes what can best be described as "exemplar-based agglomerative clustering", i.e. agglomerative clustering whose merging objective is geared towards the identification of meaningful exemplars. Analogously to apcluster(), aggExCluster() supports two variants — it can either be called for a similarity measure and data or on matrix of pairwise similarities.

#### 5.1 Getting started

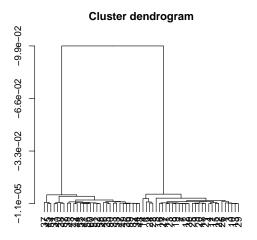
Let us start with a simple example:

```
aggres1a <- aggExCluster(negDistMat(r=2), x1)
aggres1a
##
## AggExResult object</pre>
```

```
##
## Number of samples = 60
## Maximum number of clusters = 60
```

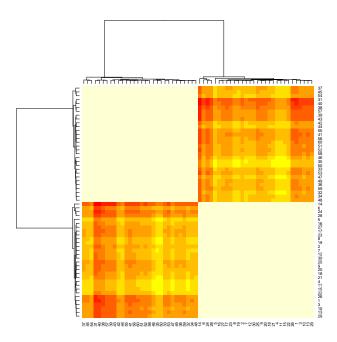
The output object aggres1a contains the complete cluster hierarchy. As obvious from the above example, the show() method only displays the most basic information. Calling plot() on an object that was the result of aggExCluster() (an object of class AggExResult), a dendrogram is plotted:

```
plot(aggres1a)
```



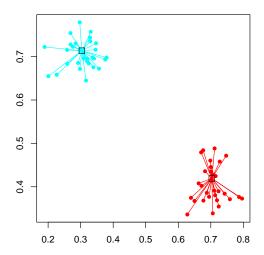
The heights of the merges in the dendrogram correspond to the merging objective: the higher the vertical bar of a merge, the less similar the two clusters have been. The dendrogram, therefore, clearly indicates two clusters. Heatmaps can be produced analogously as for APResult objects with the additional property that dendrograms are displayed on the top and on the left:

```
heatmap(aggres1a, s1)
```



Once we have confirmed the number of clusters, which is clearly 2 according to the dendrogram and the heatmap above, we can extract the level with two clusters from the cluster hierarchy. In concordance with standard R terminology, the function for doing this is called cutree():

```
cl1a <- cutree(aggres1a, k=2)</pre>
cl1a
##
## ExClust object
##
## Number of samples
## Number of clusters
##
## Exemplars:
##
      49 23
## Clusters:
##
      Cluster 1, exemplar 49:
##
         37 45 54 31 40 38 57 39 43 42 44 55 41 56 60 51 52 58 46 35 50 33
##
         53 47 49 36 59 32 34 48
      Cluster 2, exemplar 23:
##
##
         14 6 24 28 5 16 27 17 23 8 19 2 7 12 30 25 9 20 18 21 4 11 15 22
##
         26 1 3 10 13 29
plot(cl1a, x1)
```



#### 5.2 Merging clusters obtained from affinity propagation

The most important application of aggExCluster() (and the reason why it is part of the apcluster package) is that it can be used for creating a hierarchy of clusters starting from a set of clusters previously computed by affinity propagation. The examples in Section 4 indicate that it may sometimes be tricky to define the right input preference. Exemplar-based agglomerative clustering on affinity propagation results provides an additional tool for finding the right number of clusters.

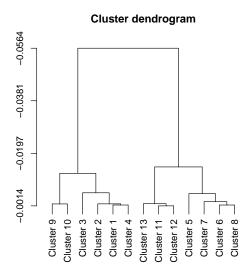
Let us revisit the four-cluster example from Section 4. We can apply aggExCluster() to an affinity propagation result if we run it on the affinity propagation result supplied as second argument x:

```
aggres2a <- aggExCluster(x=apres2c)
aggres2a

##
## AggExResult object
##
## Number of samples = 105
## Maximum number of clusters = 13</pre>
```

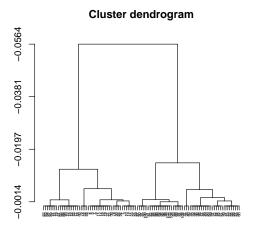
The result apres2c had 13 clusters. aggExCluster() successively joins these clusters until only one cluster is left. The dendrogram of this cluster hierarchy is given as follows:

```
plot(aggres2a)
```



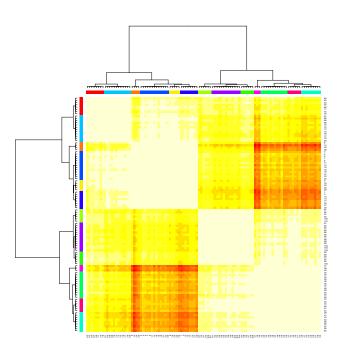
If one wants to see the original samples in the dendrogram of the cluster hierarchy, the showSamples=TRUE option can be used. In this case, it is recommended to reduce the font size of the labels via the nodePar parameter (see ?plot.dendrogram and the examples therein):





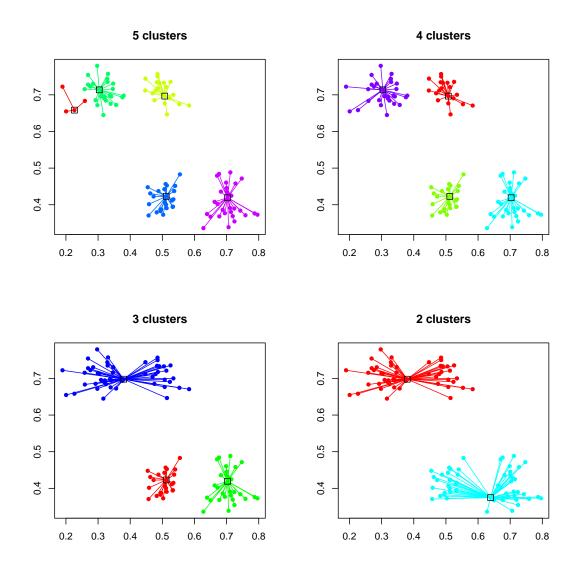
The following heatmap coincides with the one shown in Section 4 above. This is not surprising, since the heatmap plot for an affinity propagation result uses aggExCluster() internally to arrange the clusters:

heatmap(aggres2a)



Once we are more or less sure about the number of clusters, we extract the right clustering level from the hierarchy. For demonstation purposes, we do this for  $k=5,\ldots,2$  in the following plots:

```
par(mfrow=c(2,2))
for (k in 5:2)
    plot(aggres2a, x2, k=k, main=paste(k, "clusters"))
```



There is one obvious, but important, condition: applying aggExcluster() to an affinity propagation result only makes sense if the number of clusters to start from is at least as large as the number of true clusters in the data set. Clearly, if the number of clusters is already too small, then merging will make the situation only worse.

#### 5.3 Details on the merging objective

Like any other agglomerative clustering method (see, e.g., [6, 10, 13]), aggExCluster() merges clusters until only one cluster containing all samples is obtained. In each step, two clusters are merged into one, i.e. the number of clusters is reduced by one. The only aspect in which aggExCluster() differs from other methods is the merging objective.

Suppose we consider two clusters for possible merging, each of which is given by an index

set:

$$I = \{i_1, \dots, i_{n_I}\}$$
 and  $J = \{j_1, \dots, j_{n_J}\}$ 

Then we first determine the potential *joint exemplar* ex(I, J) as the sample that maximizes the average similarity to all samples in the joint cluster  $I \cup J$ :

$$\operatorname{ex}(I,J) = \operatorname*{argmax}_{i \in I \cup J} \frac{1}{n_I + n_J} \cdot \sum_{j \in I \cup J} S_{ij}$$

Recall that S denotes the similarity matrix and  $S_{ij}$  corresponds to the similarity of the *i*-th and the *j*-th sample. Then the merging objective is computed as

$$\operatorname{obj}(I,J) = \frac{1}{2} \cdot \left( \frac{1}{n_I} \cdot \sum_{j \in I} S_{\operatorname{ex}(I,J)j} + \frac{1}{n_J} \cdot \sum_{k \in J} S_{\operatorname{ex}(I,J)k} \right),$$

which can be best described as "balanced average similarity to the joint exemplar". In each step, aggExCluster() considers all pairs of clusters in the current cluster set and joins that pair of clusters whose merging objective is maximal. The rationale behind the merging objective is that those two clusters should be joined that are best described by a joint exemplar.

### 6 Leveraged Affinity Propagation

Leveraged affinity propagation is based on the idea that, for large data sets with many samples, the cluster structure is already visible on a subset of the samples. Instead of evaluating the similarity matrix for all sample pairs, the similarities of all samples to a subset of samples are computed — resulting in a non-square similarity matrix. Clustering is performed on this reduced similarity matrix allowing for clustering large data sets more efficiently.

In this form of clustering, several rounds of affinity propagation are executed with different sample subsets — iteratively improving the clustering result. The implementation is based on the Matlab code of Frey and Dueck provided on the AP Web page<sup>1</sup>. Apart from dynamic improvements through reduced amount of distance calculations and faster clustering, the memory consumption is also reduced not only in terms of the memory used for storing the similarity matrix, but also in terms of memory used by the clustering algorithm internally.

The two main parameters controlling leveraged AP clustering are the fraction of data points that should be selected for clustering (parameter frac) and the number of sweeps or repetitions of individual clustering runs (parameter sweeps). Initially, a sample subset is selected randomly. For the subsequent repetitions, the exemplars of the previous run are kept in the sample subset and the other samples in the subset are chosen randomly again. The best result of all sweeps with the highest net similarity is kept as final clustering result.

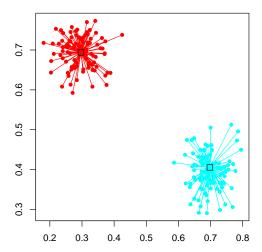
When called with a similarity measure and a dataset the function apclusterL() performs both the calculation of similarities and leveraged affinity propagation. In the example below, we use 10% of the samples and run 5 repetitions. The function implementing the similarity measure can either be passed as a function or as a function name (which must of course be resolvable in the current environment). Additional parameters for the distance calculation can be passed to apclusterL() which passes them on to the function implementing the similarity measure via the . . . argument list. In any case, this function must be implemented such that it expects the data in

its first argument x (a subsettable data structure, such as, a vector, matrix, data frame, or list) and that it takes the selection of "column objects" as a second argument sel which must be a set of column indices. The functions negDistMat(), expSimMat(), linSimMat(), corSimMat(), and linKernel() provided by the apcluster package also support the easy creation of parameter-free similarity measures (in R terminology called "closures"). We recommend this variant, as it is safer in terms of possible name conflicts between arguments of apclusterL() and arguments of the similarity function.

Here is an example that makes use of a closure for defining the similarity measure:

```
cl5 <- cbind(rnorm(100, 0.3, 0.05), rnorm(100, 0.7, 0.04))
cl6 \leftarrow cbind(rnorm(100, 0.70, 0.04), rnorm(100, 0.4, 0.05))
x3 <- rbind(c15, c16)
apres3 <- apclusterL(s=negDistMat(r=2), x=x3, frac=0.1, sweeps=5, p=-0.2)
apres3
##
## APResult object
##
## Number of samples
                            200
## Number of sel samples =
                            20
                                   (10\%)
## Number of sweeps
                            5
## Number of iterations = 124
## Input preference
                         = -0.2
## Sum of similarities = -0.7671381
## Sum of preferences
                            -0.4
## Net similarity
                         = -1.167138
## Number of clusters
                            2
##
## Exemplars:
##
      26 124
## Clusters:
##
      Cluster 1, exemplar 26:
         1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
         26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
##
         48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69
##
         70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91
##
         92 93 94 95 96 97 98 99 100
##
##
      Cluster 2, exemplar 124:
##
         101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116
         117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132
##
##
         133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148
##
         149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164
         165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180
##
         181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196
##
##
        197 198 199 200
```

```
plot(apres3, x3)
```



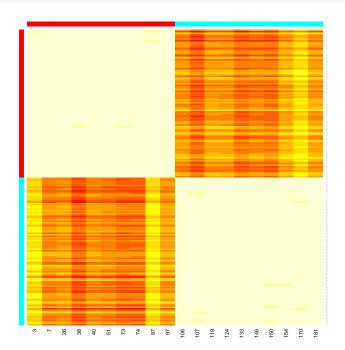
The function apclusterL() creates a result object of S4 class APResult that contains the same information as for standard AP. Additionally, the selected sample subset, the associated rectangular similarity matrix for the best sweep (provided that includeSim=TRUE) and the net similarities of all sweeps are returned in this object.

```
dim(apres3@sim)
## [1] 200
            20
apres3@sel
##
     3
            26
                38
                             73
                                  74
                                      87
                                          97 106 107 119 124 133 146 150 154
                     40
                         51
                                          97 106 107 119 124 133 146 150 154
##
     3
         7
            26
                38
                     40
                         51
                             73
                                  74
                                      87
## 170 181
## 170 181
apres3@netsimLev
## [1] -1.248064 -1.211001 -1.179808 -1.167138 -1.167138
```

The result returned by leveraged affinity propagation can be used for further processing in the same way as a result object returned from apcluster(), e.g., merging of clusters with agglomerative clustering can be performed.

For heatmap plotting either the parameter includeSim=TRUE must be set in apcluster() or apclusterL() to make the similarity matrix available in the result object or the similarity matrix must be passed as second parameter to heatmap() explicitly. The heatmap for leveraged AP looks slightly different compared to the heatmap for affinity propagation because the number of samples is different in both dimensions.

#### heatmap(apres3)



Often selected samples will be chosen as exemplars because, only for them, the full similarity information is available. This means that the fraction of samples should be selected in a way such that a considerable number of samples is available for each expected cluster. Please also note that a data set of the size used in this example can easily be clustered with regular affinity propagation. The data set was kept small to keep the package build time short and the amount of data output in the manual reasonable.

For users requiring a higher degree of flexibility, e.g., for a customized selection of the sample subset, apclusterL() called with a rectangular similiarity matrix performs affinity propagation on a rectangular similarity matrix. See the source code of apclusterL() with signature s=function and x=ANY for an example how to embed apclusterL() into a complete loop performing leveraged AP. The package-provided functions for distance calculation support the generation of rectangular similarity matrices (see Chapter 9).

# 7 Sparse Affinity Propagation

Starting with Version 1.4.0 of the apcluster package, the functions apcluster(), apclusterK(), and preferenceRange() can also handle similarity matrices as defined by the Matrix package.

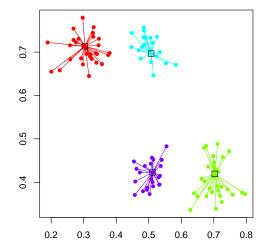
While all dense matrix formats are converted to standard R matrices, sparse matrices are converted internally to dgTMatrix objects. For these sparse matrices, special implementations of the apcluster(), apclusterK(), and preferenceRange() are available that fully exploit the sparseness of the matrices and may require much less operations if the matrix is sufficiently sparse.

In order to demonstrate that, consider the following example:

```
dsim <- negDistMat(x2, r=2)</pre>
ssim <- as.SparseSimilarityMatrix(dsim, lower=-0.2)</pre>
str(ssim)
## Formal class 'dgTMatrix' [package "Matrix"] with 6 slots
##
                : int [1:9296] 1 2 3 4 5 6 7 8 9 10 ...
                 : int [1:9296] 0 0 0 0 0 0 0 0 0 ...
##
##
     ..@ Dim
                : int [1:2] 105 105
     .. @ Dimnames:List of 2
##
##
     .. ..$ : NULL
##
     ...$ : NULL
##
                  : num [1:9296] -0.00506 -0.00162 -0.01553 -0.01426 -0.00513
     ..@ factors : list()
```

The function as .SparseSimilarityMatrix() converts the dense similarity matrix dsim into a sparse similarity matrix by removing all pairwise similarities that are -0.2 or lower. Note that this is only for demonstration purposes. If the size of data permits that, it is advisable to use the entire dense similarity matrix. Anyway, let us run sparse AP on this similarity matrix:

```
sapres <- apcluster(ssim, q=0)
plot(sapres, x2)</pre>
```

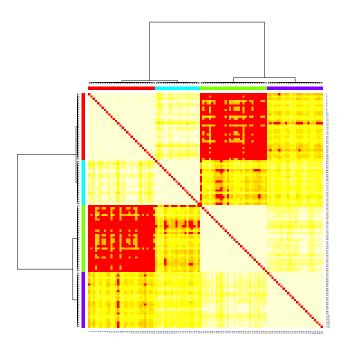


The functions preferenceRange() and apclusterK() work in the same way as for dense similarity matrices:

```
preferenceRange(ssim)
## [1] -3.718130e+00 -1.596955e-06
apclusterK(ssim, K=2)
## Trying p = -0.003719726
     Number of clusters: 19
## Trying p = -0.03718288
##
     Number of clusters: 5
## Trying p = -0.3718145
     Number of clusters: 4
##
## Trying p = -1.859066 (bisection step no. 1)
##
     Number of clusters: 2
##
## Number of clusters: 2 for p = -1.859066
##
## APResult object
##
## Number of samples
                       = 105
## Number of iterations = 128
## Input preference =
                           -1.859066
## Sum of similarities = -1.389361
## Sum of preferences = -3.718132
## Net similarity
                       = -5.107492
## Number of clusters
                       = 2
##
## Exemplars:
##
     4 54
## Clusters:
     Cluster 1, exemplar 4:
##
##
        1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
        26 27 28 29 30 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77
##
##
        78 79 80
##
     Cluster 2, exemplar 54:
##
        31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52
        53 54 55 56 57 58 59 60 81 82 83 84 85 86 87 88 89 90 91 92 93 94
##
##
        95 96 97 98 99 100 101 102 103 104 105
```

The functions aggExCluster() and heatmap() have been extended to be able to handle sparse matrices. Note, however, that these functions are not yet exploiting sparsity properly. Instead, they convert all inputs to dense matrices before processing them, which may lead to memory and/or performance issues for large data sets.

heatmap(sapres, ssim)



The above heatmap illustrates that values that are not stored in the sparse similarity matrix are filled up with low values (see the red areas between some pairs of samples that belong to different clusters). Actually, each missing value is replaced with

$$\min(s) - (\max(s) - \min(s)) = 2 \cdot \min(s) - \max(s),$$

where  $\min(s)$  and  $\max(s)$  denote the smallest and the largest similarity value specified in the sparse similarity matrix s, respectively. The same replacement takes place when aggExCluster() converts sparse similarity matrices to dense ones.

## 8 A Toy Example with Biological Sequences

As noted in the introduction above, one of the goals of this package is to leverage affinity propagation in bioinformatics applications. In order to demonstrate the usage of the package in a biological application, we consider a small toy example here.

The package comes with a toy data set ch22Promoters that consists of sub-sequences of promoter regions of 150 random genes from the human chromosome no. 22 (according to the human genome assembly hg18). Each sequence consists of the 1000 bases upstream of the transcription start site of each gene. Suppose we want to cluster these sequences in order to find out whether groups of promoters can be identified on the basis of the sequence only and, if so, to identify exemplars that are most typical for these groups.

```
library(Biostrings)
filepath <- system.file("examples", "ch22Promoters.fasta",</pre>
                       package="apcluster")
ch22Promoters <- readDNAStringSet(filepath)</pre>
ch22Promoters
##
    A DNAStringSet instance of length 150
##
       width seq
                                                       names
     [1] 1000 AGACTTAAGGGACCTGGT...CGTGTGCGCATGCGCAGC NM_001169111
##
     [2] 1000 CCCGGCTAATTTTTTTGT...CGGAGTCCGGGCGAGGTG NM_012324
##
##
     [3] 1000 CACATGTGCCCTCTGGGC...CAAACGCAGCGCCAGACA NM_144704
##
     [4] 1000 AAGCATGGTGGGATTGGC...GAATGGTCCCGCGGCTCC NM_002473
##
     [5] 1000 TTTAGAGAACTGGGTCTT...CAGAGCCCAGCGGGAGCG NM_001184970
##
     ## [146] 1000 TCCGCCTCCTGGGTTCAA...TCGGGGAGGGGCAGTAAG NM_032608
## [147] 1000 ACTAAACTTAGTATATTA...CGCGGGTGGGCGGGCCT NM_003560
## [148] 1000 AGGATCACATCAGCTAAC...TTCAACCCTAGATCCACC NM_001166242
## [149] 1000 AGGCAGGAGAATCGATTG...GGATCATGAGGTCAGGAG NM_001165877
## [150] 1000 GCTCAGGATGAAATAGGC...GAGCTGCTGGGAAGTTGT NM_145343
```

Obviously, these are classical nucleotide sequences, each of which is identified by the RefSeq identifier of the gene the promoter sequence stems from.

In order to compute a similarity matrix for this data set, we choose (without further justification, just for demonstration purposes) the simple *spectrum kernel* [8] with a sub-sequence length of k=6. We use the implementation from the kebabs package [11] to compute the similarity matrix in a convenient way:

Now we run affinity propagation on the similarity matrix:

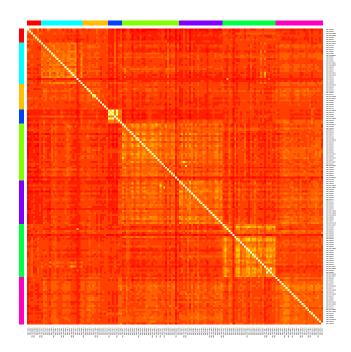
```
promAP <- apcluster(promSim, q=0)
promAP
##</pre>
```

```
## APResult object
##
## Number of samples = 150
## Number of iterations = 187
## Input preference = 0.01367984
## Sum of similarities = 57.79686
## Sum of preferences = 0.1094387
## Net similarity
                        = 57.9063
## Number of clusters = 8
##
## Exemplars:
      NM_001199580 NM_022141 NM_152868 NM_001128633 NM_052945 NM_001099294
##
##
      NM_152513 NM_080764
## Clusters:
##
      Cluster 1, exemplar NM_001199580:
##
         NM_001169111 NM_012324 NM_144704 NM_002473 NM_005198 NM_004737
         NM_007194 NM_014292 NM_001199580 NM_032758 NM_003325 NM_014876
##
        NM_000407 NM_053004 NM_023004 NM_001130517 NM_002882 NM_001169110
##
        NM_020831 NM_001195071 NM_031937 NM_001164502 NM_152299 NM_014509
##
##
        NM_138433 NM_006487 NM_005984 NM_001085427 NM_013236
##
      Cluster 2, exemplar NM_022141:
         NM_001184970 NM_002883 NM_001051 NM_014460 NM_153615 NM_022141
##
##
         NM_001137606 NM_001184971 NM_053006 NM_015367 NM_000395 NM_012143
##
         NM 004147
##
      Cluster 3, exemplar NM_152868:
##
         {\tt NM\_015140\ NM\_138415\ NM\_001195072\ NM\_001164501\ NM\_014339\ NM\_152868}
         NM_033257 NM_014941 NM_033386 NM_007311 NM_017911 NM_007098
##
##
         NM_001670 NM_015653 NM_032775 NM_001172688 NM_001136029
##
         NM_001024939 NM_002969 NM_052906 NM_152511 NM_001001479
##
      Cluster 4, exemplar NM_001128633:
         NM_001128633 NM_002133 NM_015672 NM_013378 NM_001128633 NM_000106
##
##
         NM 152855
      Cluster 5, exemplar NM_052945:
##
##
         NM_001102371 NM_017829 NM_020243 NM_002305 NM_030882 NM_000496
         NM_022720 NM_024053 NM_052945 NM_018006 NM_014433 NM_032561
##
##
         NM_001142964 NM_181492 NM_003073 NM_015366 NM_005877 NM_148674
##
         NM_005008 NM_017414 NM_000185 NM_001135911 NM_001199562 NM_003935
##
         NM_003560 NM_001166242 NM_001165877
      Cluster 6, exemplar NM_001099294:
##
##
         NM 004900 NM 001097 NM 174975 NM 004377 NM 001099294 NM 004121
##
         NM_001146288 NM_002415 NM_001159546 NM_004327 NM_152426 NM_004861
         NM_001193414 NM_001145398 NM_015715 NM_021974 NM_001159554
##
##
         {\tt NM\_001171660\ NM\_015124\ NM\_006932\ NM\_152906\ NM\_001002034\ NM\_000754}
##
        NM 145343
      Cluster 7, exemplar NM_152513:
##
##
         NM_001135772 NM_007128 NM_014227 NM_203377 NM_152267 NM_017590
```

```
## NM_001098270 NM_138481 NM_012401 NM_014303 NM_001144931
## NM_001098535 NM_000262 NM_152510 NM_003347 NM_001039141
## NM_001014440 NM_152513 NM_015330 NM_138338 NM_032608
## Cluster 8, exemplar NM_080764:
## NM_175878 NM_003490 NM_001039366 NM_001010859 NM_014306 NM_080764
## NM_001164104
```

So we obtain 8 clusters in total. The corresponding heatmap looks as follows (note that we have to supply the similarity matrix, as it is not included by default if apcluster() is called with the similarity matrix; the reason is that, for large data sets, it is more memory-efficient not to make multiple copies of the similarity matrix; see also 10.1):

```
heatmap(promAP, promSim, Rowv=FALSE, Colv=FALSE)
```

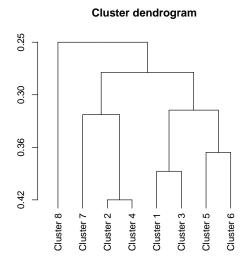


Let us now run agglomerative clustering to further join clusters.

```
promAgg <- aggExCluster(promSim, promAP)</pre>
```

The resulting dendrogram is given as follows:

```
plot(promAgg)
```



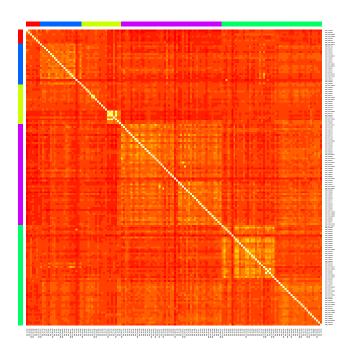
The dendrogram does not give a very clear indication about the best number of clusters. Let us adopt the viewpoint for a moment that 5 clusters are reasonable.

```
prom5 <- cutree(promAgg, k=5)</pre>
prom5
##
## ExClust object
##
## Number of samples
## Number of clusters =
##
## Exemplars:
      NM_152513 NM_080764 NM_001128633 NM_001199580 NM_052945
##
## Clusters:
##
      Cluster 1, exemplar NM_152513:
         NM_001135772 NM_007128 NM_014227 NM_203377 NM_152267 NM_017590
##
         NM_001098270 NM_138481 NM_012401 NM_014303 NM_001144931
##
         NM_001098535 NM_000262 NM_152510 NM_003347 NM_001039141
##
##
         NM_001014440 NM_152513 NM_015330 NM_138338 NM_032608
##
      Cluster 2, exemplar NM_080764:
##
         NM_175878 NM_003490 NM_001039366 NM_001010859 NM_014306 NM_080764
         NM_001164104
##
##
      Cluster 3, exemplar NM_001128633:
         NM_001184970 NM_002883 NM_001051 NM_014460 NM_153615 NM_022141
##
##
         NM_001137606 NM_001184971 NM_053006 NM_015367 NM_000395 NM_012143
##
         NM_004147 NM_001128633 NM_002133 NM_015672 NM_013378 NM_001128633
##
         NM_000106 NM_152855
      Cluster 4, exemplar NM_001199580:
##
```

```
NM_001169111 NM_012324 NM_144704 NM_002473 NM_005198 NM_004737
##
##
         NM_007194 NM_014292 NM_001199580 NM_032758 NM_003325 NM_014876
         NM_000407 NM_053004 NM_023004 NM_001130517 NM_002882 NM_001169110
##
         NM_020831 NM_001195071 NM_031937 NM_001164502 NM_152299 NM_014509
##
##
         NM_138433 NM_006487 NM_005984 NM_001085427 NM_013236 NM_015140
         NM_138415 NM_001195072 NM_001164501 NM_014339 NM_152868 NM_033257
##
         NM_014941 NM_033386 NM_007311 NM_017911 NM_007098 NM_001670
##
         NM_015653 NM_032775 NM_001172688 NM_001136029 NM_001024939
##
         NM_002969 NM_052906 NM_152511 NM_001001479
##
      Cluster 5, exemplar NM_052945:
##
##
         NM_001102371 NM_017829 NM_020243 NM_002305 NM_030882 NM_000496
         NM_022720 NM_024053 NM_052945 NM_018006 NM_014433 NM_032561
##
##
         NM_001142964 NM_181492 NM_003073 NM_015366 NM_005877 NM_148674
         NM_005008 NM_017414 NM_000185 NM_001135911 NM_001199562 NM_003935
##
##
         NM_003560 NM_001166242 NM_001165877 NM_004900 NM_001097 NM_174975
##
         NM_004377 NM_001099294 NM_004121 NM_001146288 NM_002415
         NM_001159546 NM_004327 NM_152426 NM_004861 NM_001193414
##
##
         NM_001145398 NM_015715 NM_021974 NM_001159554 NM_001171660
         NM_015124 NM_006932 NM_152906 NM_001002034 NM_000754 NM_145343
##
```

The final heatmap looks as follows:

heatmap(prom5, promSim, Rowv=FALSE, Colv=FALSE)



# 9 Similarity Matrices

Apart from the obvious monotonicity "the higher the value, the more similar two samples", affinity propagation does not make any specific assumption about the similarity measure. Negative squared distances must be used if one wants to minimize squared errors [5]. Apart from that, the choice and implementation of the similarity measure is left to the user.

Our package offers a few more methods to obtain similarity matrices. The choice of the right one (and, consequently, the objective function the algorithm optimizes) still has to be made by the user.

All functions described in this section assume the input data matrix to be organized such that each row corresponds to one sample and each column corresponds to one feature (in line with the standard function dist). If a vector is supplied instead of a matrix, each single entry is interpreted as a (one-dimensional) sample.

# **9.1** The function negDistMat()

The function negDistMat(), in line with Frey and Dueck, allows for computing negative distances for a given set of real-valued data samples. If called with the first argument x, a similarity matrix with pairwise negative distances is returned:

The function negDistMat() provides the same set of distance measures and parameters as the standard function dist() (except for method="binary" which makes little sense for real-valued data). Presently, negDistMat() provides the following variants of computing the distance  $d(\mathbf{x}, \mathbf{y})$  of two data samples  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$ :

#### **Euclidean:**

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

use method="euclidean" or do not specify argument method (since this is the default);

## **Maximum:**

$$d(\mathbf{x}, \mathbf{y}) = \max_{i=1}^{n} |x_i - y_i|$$

use method="maximum";

## Sum of absolute distances / Manhattan:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{n} |x_i - y_i|$$

use method="manhattan";

#### Canberra:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{n} \frac{|x_i - y_i|}{|x_i + y_i|}$$

summands with zero denominators are not taken into account; use method="canberra";

### Minkowski:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^{n} (x_i - y_i)^p\right)^{\frac{1}{p}}$$

use method="minkowski" and specify p using the additional argument p (default is p=2, resulting in the standard Euclidean distance);

The function negDistMat() then takes the distances computed with one of the variants listed above and returns -1 times the r-th power of it, i.e.,

$$s(\mathbf{x}, \mathbf{y}) = -d(\mathbf{x}, \mathbf{y})^r. \tag{1}$$

The exponent r can be adjusted with the argument r. The default is r=1, hence, one has to supply r=2 to obtain negative squared distances as in the examples in previous sections.

Here are some examples:

Standard Euclidean distance:

```
megDistMat(ex)

## 1 2 3 4 5

## 1 0.0000000 -1.2688578 -0.8831761 -0.7071068 -1.2041595

## 2 -1.2688578 0.0000000 -0.8660254 -1.4177447 -0.8000000

## 3 -0.8831761 -0.8660254 0.0000000 -0.6164414 -0.5196152

## 4 -0.7071068 -1.4177447 -0.6164414 0.0000000 -1.0246951

## 5 -1.2041595 -0.8000000 -0.5196152 -1.0246951 0.0000000
```

Squared Euclidean distance:

```
negDistMat(ex, r=2)

## 1 2 3 4 5

## 1 0.00 -1.61 -0.78 -0.50 -1.45

## 2 -1.61 0.00 -0.75 -2.01 -0.64

## 3 -0.78 -0.75 0.00 -0.38 -0.27

## 4 -0.50 -2.01 -0.38 0.00 -1.05

## 5 -1.45 -0.64 -0.27 -1.05 0.00
```

Maximum norm-based distance:

```
negDistMat(ex, method="maximum")

## 1 2 3 4 5

## 1 0.0 -1.0 -0.7 -0.5 -1.0

## 2 -1.0 0.0 -0.7 -1.0 -0.8

## 3 -0.7 -0.7 0.0 -0.5 -0.5

## 4 -0.5 -1.0 -0.5 0.0 -1.0

## 5 -1.0 -0.8 -0.5 -1.0 0.0
```

Sum of absolute distances (aka Manhattan distance):

```
negDistMat(ex, method="manhattan")

## 1 2 3 4 5

## 1 0.0 -2.1 -1.4 -1.0 -1.9

## 2 -2.1 0.0 -1.3 -2.1 -0.8

## 3 -1.4 -1.3 0.0 -1.0 -0.7

## 4 -1.0 -2.1 -1.0 0.0 -1.3

## 5 -1.9 -0.8 -0.7 -1.3 0.0
```

Canberra distance:

```
negDistMat(ex, method="canberra")

## 1 2 3 4 5

## 1 0.000000 -2.600000 -1.9444444 -1.181818 -1.8307692

## 2 -2.600000 0.000000 -1.66666667 -2.200000 -1.0000000

## 3 -1.944444 -1.666667 0.0000000 -1.676471 -0.7333333

## 4 -1.181818 -2.200000 -1.6764706 0.0000000 -1.3111111

## 5 -1.830769 -1.000000 -0.7333333 -1.311111 0.0000000
```

Minkowski distance for p = 3 (3-norm):

If called without the data argument x, a function object is returned that can be supplied to clustering functions — as in the majority of the above examples:

```
sim <- negDistMat(r=2)</pre>
is.function(sim)
## [1] TRUE
apcluster(sim, x1)
##
## APResult object
##
## Number of samples = 60
## Number of iterations = 130
## Input preference = -0.1582302
## Sum of similarities = -0.1974598
## Sum of preferences = -0.3164603
## Net similarity = -0.5139202
## Number of clusters
                       = 2
##
## Exemplars:
     23 49
##
## Clusters:
     Cluster 1, exemplar 23:
##
        1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
##
        26 27 28 29 30
##
     Cluster 2, exemplar 49:
        31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52
##
##
        53 54 55 56 57 58 59 60
```

# 9.2 Other similarity measures

The package apcluster offers four more functions for creating similarity matrices for real-valued data:

Exponential transformation of distances: the function expSimMat() works in the same way as the function negDistMat(). The difference is that, instead of the transformation (1), it uses the following transformation:

$$s(\mathbf{x}, \mathbf{y}) = \exp\left(-\left(\frac{d(\mathbf{x}, \mathbf{y})}{w}\right)^r\right)$$

Here the default is r=2. It is clear that r=2 in conjunction with method="euclidean" results in the well-known *Gaussian kernel | RBF kernel* [4, 9, 12], whereas r=1 in conjunction with method="euclidean" results in the similarity measure that is sometimes called *Laplace kernel* [4, 9]. Both variants (for non-Euclidean distances as well) can also be interpreted as *fuzzy equality/similarity relations* [2].

**Linear scaling of distances with truncation:** the function linSimMat() uses the transformation

$$s(\mathbf{x}, \mathbf{y}) = \max\left(1 - \frac{d(\mathbf{x}, \mathbf{y})}{w}, 0\right)$$

which is also often interpreted as a fuzzy equality/similarity relation [2].

Correlation: the function corSimMat() interprets the rows of its argument x (matrix or data frame) as multivariate observations and computes similarities as pairwise correlations. The function corSimMat() is actually a wrapper around the standard function cor(). Consequently, the method argument allows for selecting the type of correlation to compute (Pearson, Spearman, or Kendall).

**Linear kernel:** scalar products can also be interpreted as similarity measures, a view that is often adopted by kernel methods in machine learning. In order to provide the user with this option as well, the function linKernel() is available. For two data samples  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$ , it computes the similarity as

$$s(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{n} x_i \cdot y_i.$$

The function has one additional argument, normalize (default: FALSE). If normalize is set to TRUE, values are normalized to the range [-1, +1] in the following way:

$$s(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^{n} x_i \cdot y_i}{\sqrt{\left(\sum_{i=1}^{n} x_i^2\right) \cdot \left(\sum_{i=1}^{n} y_i^2\right)}}$$

Entries for which at least one of the two factors in the denominator is zero are set to zero (however, the user should be aware that this should be avoided anyway).

For the same example data as above, we obtain the following for the RBF kernel:

```
## 1 2 3 4 5

## 1 1.0000000 0.1998876 0.4584060 0.6065307 0.2345703

## 2 0.1998876 1.0000000 0.4723666 0.1339887 0.5272924

## 3 0.4584060 0.4723666 1.0000000 0.6838614 0.7633795

## 4 0.6065307 0.1339887 0.6838614 1.0000000 0.3499377

## 5 0.2345703 0.5272924 0.7633795 0.3499377 1.0000000
```

## Laplace kernel:

```
## 1 2 3 4 5
## 1 1.0000000 0.2811526 0.4134676 0.4930687 0.2999440
## 2 0.2811526 1.0000000 0.4206200 0.2422598 0.4493290
## 3 0.4134676 0.4206200 1.0000000 0.5398622 0.5947493
## 4 0.4930687 0.2422598 0.5398622 1.0000000 0.3589059
## 5 0.2999440 0.4493290 0.5947493 0.3589059 1.0000000
```

#### Pearson correlation coefficient:

# Spearman rank correlation coefficient:

```
corSimMat(ex, method="spearman")

## 1 2 3 4 5

## 1 1.0 -0.5 -0.5 0.5 -1.0

## 2 -0.5 1.0 -0.5 -1.0 0.5

## 3 -0.5 -0.5 1.0 0.5 0.5

## 4 0.5 -1.0 0.5 1.0 -0.5

## 5 -1.0 0.5 0.5 -0.5 1.0
```

## Linear scaling of distances with truncation:

```
linSimMat(ex, w=1.2)
```

```
## 1 1 2 3 4 5

## 1 1.0000000 0.0000000 0.2640199 0.4107443 0.0000000

## 2 0.0000000 1.0000000 0.2783122 0.0000000 0.3333333

## 3 0.2640199 0.2783122 1.0000000 0.4862988 0.5669873

## 4 0.4107443 0.0000000 0.4862988 1.0000000 0.1460874

## 5 0.0000000 0.3333333 0.5669873 0.1460874 1.0000000
```

### Linear kernel:

```
linKernel(ex[2:5,])

##     1     2     3     4

## 1 1.04 0.52 0.06 1.04

## 2 0.52 0.75 0.73 1.08

## 3 0.06 0.73 1.09 0.86

## 4 1.04 1.08 0.86 1.68
```

### Normalized linear kernel:

All of these functions work in the same way as negDistMat(): if called with argument x, a similarity matrix is returned, otherwise a function is returned.

## 9.3 Rectangular similarity matrices

With the introduction of leveraged affinity propagation, distance calculations are entirely performed within the apcluster package. The code is based on a customized version of the dist() function from the stats package. In the following example, a rectangular similarity matrix of all samples against a subset of the samples is computed:

```
sel <- sort(sample(1:nrow(x1), ceiling(0.08 * nrow(x1))))
sel

## [1] 4 9 13 14 56

s1r <- negDistMat(x1, sel, r=2)
dim(s1r)</pre>
```

```
## [1] 60 5

s1r[1:7,]

## 4 9 13 14 56

## 1 -0.015532442 -0.003720467 -0.0007006572 -0.007335482 -0.3375721

## 2 -0.004985640 -0.002699246 -0.0026429582 -0.015075926 -0.2628843

## 3 -0.015085144 -0.005792112 -0.0002679638 -0.004800287 -0.3194392

## 4 0.000000000 -0.004728567 -0.0131842627 -0.036886156 -0.2120570

## 5 -0.006807057 -0.009969146 -0.0091423199 -0.022165335 -0.2274750

## 6 -0.014798343 -0.008546240 -0.0020775500 -0.006370294 -0.2982598

## 7 -0.003253827 -0.003592281 -0.0049782050 -0.019700651 -0.2436157
```

The rows correspond to all samples, the columns to the sample subset. The sel parameter specifies the sample indices of the selected samples in increasing order. Rectangular similarity calculation is provided in all distance functions of the package. If the parameter sel is not specified, the quadratic similarity matrix of all sample pairs is computed.

## 9.4 Defining a custom similarity measure for leveraged affinity propagation

As mentioned in Section 6 above, leveraged affinity propagation requires the definition of a similarity measure that is supplied as a function or function name to apclusterL(). For vectorial data, the similarity measures supplied with the package (see above) may be sufficient. If other similarity measures are necessary or if the data are not vectorial, the user must supply his/her own similarity measure. The user can supply any function as argument s to apcluster(), apclusterK(), or apclusterL(), but the following rules must be obeyed in order to avoid errors and to ensure meaningful results:

- 1. The data must be supplied as first argument, which must be named x.
- 2. The second argument must be named sel and must be interpreted as a vector of indices that select a subset of data items in x.
- 3. The function must return a numeric matrix with similarities. If sel=NA, the format of the matrix must be  $length(x) \times length(x)$ . If sel is not NA, but contains indices selecting a subset, the format of the returned similarity matrix must be  $length(x) \times length(sel)$ .
- 4. Although this is not a must, it is recommended to properly set row and column names in the returned similarity matrix.

As an example, let us revisit the sequence clustering example presented in Section 8. Let us first define the function that implements the similarity measure:

```
spectrumK6 <- function(x, sel=NA)
{
   if (any(is.na(sel)))</pre>
```

```
s <- getKernelMatrix(specK6, x)
else
    s <- getKernelMatrix(specK6, x, x[sel])
as(s, "matrix")
}</pre>
```

Now we run leveraged affinity propagation on the ch22Promoters data set using this similarity measure

```
promAPL <- apclusterL(s=spectrumK6, ch22Promoters, frac=0.1, sweeps=10,</pre>
                      p=promAP@p)
promAPL
##
## APResult object
##
## Number of samples = 150
## Number of sel samples = 15
                                 (10%)
## Number of sweeps = 10
## Number of iterations = 135
## Input preference = 0.01367984
## Sum of similarities = 55.7574
## Sum of preferences = 0.1367984
## Net similarity = 55.8942
## Number of clusters = 10
##
## Exemplars:
      NM_004737 NM_017829 NM_175878 NM_020243 NM_152868 NM_001128633
##
     NM_003347 NM_001039141 NM_001001479 NM_003560
##
## Clusters:
     Cluster 1, exemplar NM_004737:
##
##
        NM_001169111 NM_012324 NM_005198 NM_004737 NM_007194 NM_014292
         NM_001199580 NM_003325 NM_014876 NM_000407 NM_053004 NM_023004
##
         NM_001130517 NM_001169110 NM_001195071 NM_001164502 NM_152299
##
##
         NM_138433 NM_006487 NM_005984 NM_001085427
##
      Cluster 2, exemplar NM_017829:
##
         NM_001102371 NM_017829 NM_030882 NM_024053 NM_052945 NM_018006
         NM_014433 NM_032561 NM_001142964 NM_181492 NM_003073 NM_015366
##
        NM_148674 NM_005008 NM_017414 NM_000185 NM_001166242 NM_001165877
##
##
      Cluster 3, exemplar NM_175878:
         NM_002883 NM_175878 NM_003490 NM_001184971 NM_001010859 NM_012143
##
##
        NM_014306 NM_080764 NM_004147
##
      Cluster 4, exemplar NM_020243:
##
         NM_020243 NM_015330 NM_001164104
      Cluster 5, exemplar NM_152868:
##
```

```
NM_015140 NM_138415 NM_001164501 NM_014339 NM_152868 NM_033257
##
##
         NM_014941 NM_033386 NM_017911 NM_007098 NM_032775 NM_001172688
##
         NM_001136029 NM_052906
      Cluster 6, exemplar NM_001128633:
##
##
         NM_001128633 NM_002133 NM_015672 NM_013378 NM_001128633 NM_000106
         NM 152855
##
      Cluster 7, exemplar NM_003347:
##
         NM_001051 NM_007128 NM_014227 NM_203377 NM_014460 NM_017590
##
         NM_001098270 NM_153615 NM_022141 NM_001137606 NM_000496 NM_053006
##
##
         NM 138481 NM 015367 NM 014303 NM 001098535 NM 000262 NM 003347
##
         NM_152513 NM_021974 NM_138338 NM_013236
      Cluster 8, exemplar NM_001039141:
##
##
         NM_001135772 NM_152267 NM_012401 NM_001144931 NM_152510
##
         NM_001039141 NM_001014440 NM_001135911 NM_032608
##
      Cluster 9, exemplar NM_001001479:
##
         NM_144704 NM_001184970 NM_004900 NM_001097 NM_174975 NM_001195072
         NM_032758 NM_002305 NM_004377 NM_001099294 NM_002882 NM_020831
##
         {\tt NM\_004121\ NM\_001146288\ NM\_031937\ NM\_002415\ NM\_001159546\ NM\_000395}
##
         NM_004327 NM_152426 NM_007311 NM_001670 NM_004861 NM_014509
##
         NM_001193414 NM_015653 NM_001145398 NM_001024939 NM_015715
##
         NM_002969 NM_152511 NM_001001479 NM_001159554 NM_001171660
##
##
         NM_015124 NM_006932 NM_152906 NM_001002034 NM_000754 NM_145343
##
      Cluster 10, exemplar NM_003560:
         NM_002473 NM_001039366 NM_022720 NM_005877 NM_001199562 NM_003935
##
##
         NM_003560
```

So we obtain 10 clusters in total.

# 9.5 Defining a custom similarity measure that creates a sparse similarity matrix

Since Version 1.4.0, similarity matrices may also be sparse (cf. Section 7). Correspondingly, the similarity measures passed to apcluster() and apclusterK() may also return sparse similarity matrices:

```
sparseSim <- function(x)
{
    as.SparseSimilarityMatrix(negDistMat(x, r=2), lower=-0.2)
}
sapres2 <- apcluster(sparseSim, x2, q=0)
sapres2
##
## APResult object
##</pre>
```

10 Miscellaneous 49

```
## Number of samples
                     = 105
## Number of iterations = 132
## Input preference
                     = -0.1998742
## Sum of similarities = -0.2820462
## Sum of preferences
                        = -0.7994967
## Net similarity
                        = -1.081543
## Number of clusters
                       = 4
##
## Exemplars:
##
      23 49 71 81
## Clusters:
      Cluster 1, exemplar 23:
##
##
        1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
##
        26 27 28 29 30
##
      Cluster 2, exemplar 49:
##
        31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52
##
        53 54 55 56 57 58 59 60
##
      Cluster 3, exemplar 71:
        61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80
##
      Cluster 4, exemplar 81:
##
        81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101
##
##
         102 103 104 105
str(similarity(sapres2))
## Formal class 'dgTMatrix' [package "Matrix"] with 6 slots
                : int [1:9296] 1 2 3 4 5 6 7 8 9 10 ...
##
##
     ..@ j
                : int [1:9296] 0 0 0 0 0 0 0 0 0 ...
                : int [1:2] 105 105
##
     ..@ Dim
    .. @ Dimnames:List of 2
##
    ....$ : NULL
    ....$ : NULL
##
##
                 : num [1:9296] -0.00506 -0.00162 -0.01553 -0.01426 -0.00513 ...
## ..@ factors : list()
```

Note that similarity measures passed to apclusterL() may not return sparse matrices. Instead, they must accept a sel argument and return a rectangular dense matrix (see Subsection 9.4 above).

# 10 Miscellaneous

## 10.1 Convenience vs. efficiency

In most of the above examples, we called a clustering method by supplying it with a similarity function and the data to be clustered. This is undoubtedly a convenient approach. Since the result-

50 10 Miscellaneous

ing output objects (unless the option includeSim=FALSE is supplied) even includes the similarity matrix, we can plot heatmaps and produce a cluster hierarchy on the basis of the clustering result without the need to supply the similarity matrix explicitly.

For large data sets, however, this convenient approach has some disadvantages:

- If the clustering algorithm is run several times on the same data set (e.g., for different parameters), the similarity matrix is recomputed every time.
- Every clustering result (depending on the option includeSim) usually includes a copy of the similarity matrix.

For these reasons, depending on the actual application scenario, users should consider computing the similarity matrix beforehand. This strategy, however, requires some extra effort for subsequent processing, i.e. the similarity must be supplied as an extra argument in subsequent processing.

# 10.2 Clustering named objects

The function apcluster() and all functions for computing distance matrices are implemented to recognize names of data objects and to correctly pass them through computations. The mechanism is best described with a simple example:

```
x3 <- c(1, 2, 3, 7, 8, 9)
names(x3) <- c("a", "b", "c", "d", "e", "f")
s3 <- negDistMat(x3, r=2)
```

So we see that the names attribute must be used if a vector of named one-dimensional samples is to be clustered. If the data are not one-dimensional (a matrix or data frame), object names must be stored in the row names of the data matrix.

All functions for computing similarity matrices recognize the object names. The resulting similarity matrix has the list of names both as row and column names.

```
s3
##
         b
              С
                  d
## a
      0 -1 -4 -36 -49 -64
## b -1 0 -1 -25 -36 -49
## c -4 -1
            0 -16 -25 -36
## d -36 -25 -16
                0 -1 -4
## e -49 -36 -25 -1
                        -1
## f -64 -49 -36 -4
colnames(s3)
## [1] "a" "b" "c" "d" "e" "f"
```

10 Miscellaneous 51

The function apcluster() and all related functions use column names of similarity matrices as object names. If object names are available, clustering results are by default shown by names.

```
apres3a <-apcluster(s3)</pre>
apres3a
##
## APResult object
##
## Number of samples
## Number of iterations = 124
## Input preference = -25
## Sum of similarities = -4
## Sum of preferences = -50
                   = -54
## Net similarity
## Number of clusters = 2
##
## Exemplars:
##
     bе
## Clusters:
##
     Cluster 1, exemplar b:
        a b c
##
##
     Cluster 2, exemplar e:
##
        def
apres3a@exemplars
## b e
## 2 5
apres3a@clusters
## [[1]]
## a b c
## 1 2 3
##
## [[2]]
## d e f
## 4 5 6
```

# 10.3 Computing a label vector from a clustering result

For later classification or comparisons with other clustering methods, it may be useful to compute a label vector from a clustering result. Our package provides an instance of the generic function

52 10 Miscellaneous

labels() for this task. As obvious from the following example, the argument type can be used to determine how to compute the label vector.

```
apres3a@exemplars

## b e
## 2 5

labels(apres3a, type="names")

## [1] "b" "b" "b" "e" "e" "e"

labels(apres3a, type="exemplars")

## [1] 2 2 2 5 5 5

labels(apres3a, type="enum")

## [1] 1 1 1 2 2 2
```

The first choice, "names" (default), uses names of exemplars as labels (if names are available, otherwise an error message is displayed). The second choice, "exemplars", uses indices of exemplars (enumerated as in the original data set). The third choice, "enum", uses indices of clusters (consecutively numbered as stored in the slot clusters — analogous to the standard implementation of cutree() or the clusters field of the list returned by the standard function kmeans()).

## 10.4 Customizing heatmaps

With Version 1.3.1, the implementation of heatmap plotting has changed significantly. The method now allows for many more customizations than before. Apart from changes in the argument list (see ?heatmap), the behavior of the method has changed as follows:

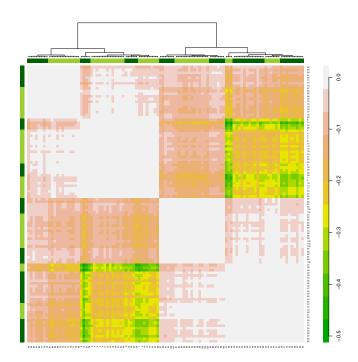
- Dendrograms are always plotted if possible. To switch off plotting of dendrograms, set Rowv and Colv to FALSE or NA. If a dendrogram should only appear to the left of the heatmap, set Colv to FALSE or NA. Analogously, set Rowv to FALSE or NA if a dendrogram should only be plotted on top of the plot (not possible if the similarity matrix is non-quadratic).
- Previously, rainbow() was used internally to determine how the bars illustrating the clusters are colored. Now users can determine the coloring of the color bars using the sideColors argument. For sideColors=NULL, a meaningful color coding is determined automatically which still uses rainbow(), but ensures that no similar colors are placed next to each other in the bar.

10 Miscellaneous 53

■ The default font sizes for displaying row/column labels have been changed to make sure that they do not overlap. This can result in quite small labels if the number of samples is larger. In any case, the user can override the sizes by making custom settings of the parameters cexRow and cexCol. Row and column labels can even be switched off entirely by setting cexRow and cexCol to 0, respectively.

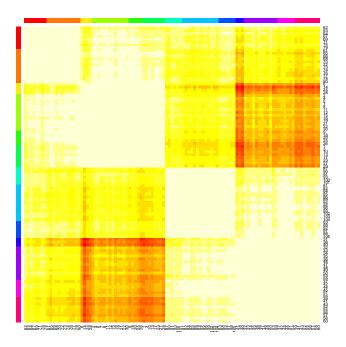
Moreover, with Version 1.4.3, the possibility to add a color legend has been integrated.

Here is an example with the vertical dendrogram switched off, an alternate color scheme, custom margins, and a color legend:



The following example reverts to the default behavior prior to Version 1.3.1: consecutive rainbow colors, no dendrograms, and traditional sizing of row/column labels:

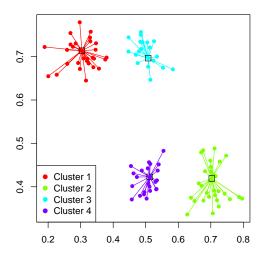
54 10 Miscellaneous



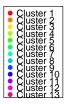
# 10.5 Adding a legend to plots of clustering results

As shown above, plot() called for an APResult object as first and a matrix or data frame as second argument plots the clustering result superimposed on a scatter plot (or a scatter plot matrix if the number of columns in the second argument exceeds 2). The clusters are shown in different colors, but it may not be clear which cluster is shown in which color. Therefore, it may be useful to show a legend along with the plot. The current implementation of plot() does not show a legend, since it is hard to determine where to actually place the legend such that no important cluster information gets occluded by the legend. Therefore, the user has to add legends manually. Actually, colors are always chosen according to a simple rule: plot() uses rainbow() to create a vector of colors that is exactly as long as the number of clusters in the APResult object. The following example shows how to plot a legend manually (with the clusters enumerated in the same way as in the APResult object):

10 Miscellaneous 55



Note that this method is only meaningful for plotting clustering results superimposed on a 2D data set. For scatter plot matrices, this does not work in a meaningful way. In such a case, the user is rather recommended to create a legend separately (in a separate graphics device/file) and to display it along with the scatter plot matrix. To create only the legend, code like the following could be used:



It still may be necessary to strip off white margins for further usage of the legend.

# 10.6 Implementation and performance issues

Prior to Version 1.2.0, apcluster() was implemented in R. Starting with version 1.2.0, the main iteration loop of apcluster() has been implemented in C++ using the Rcpp package [3], which has led to a speedup in the range of a factor or 9–10.

Note that details=TRUE requires quite an amount of additional memory. If possible, avoid this for larger data sets.

The asymptotic computational complexity of aggExCluster() is  $\mathcal{O}(l^3)$  (where l is the number of samples or clusters from which the clustering starts). This may result in excessively long computation times if aggExCluster() is used for larger data sets without using affinity propagation first. For real-world data sets, in particular, if they are large, we recommend to use affinity propagation first and then, if necessary, to use aggExCluster() to create a cluster hierarchy.

# 10.7 Using APCluster in Conjunction with KeBABS

The example in Section 8 makes use of the kebabs package [11]. The two packages have been designed to work well with each other, but since they are covering quite different domains and are available from different repositories (kebabs is in Bioconductor, while apcluster is in CRAN), we deliberately tried to avoid any hard dependencies. This trade-off between dependence and independences could not be resolved perfectly. So, it is necessary to make sure that kebabs is loaded before apcluster to make sure that, for instance, the heatmap() function from the apcluster package is not overridden by the kebabs package. In any case, we will try to eliminate these difficulties in future releases.

# 11 Special Notes for Users Upgrading from Previous Versions

# 11.1 Upgrading from a version older than 1.3.0

Version 1.3.0 has brought several fundamental changes to the architecture of the package. We tried to ensure backward compatibility with previous versions where possible. However, there are still some caveats the users should take into account:

- The functions apcluster(), apclusterK(), and aggExCluster() have been re-implemented as S4 generics, therefore, they do not have a fixed list of arguments anymore. For this reason, users are recommended to name all optional parameters.
- Heatmap plotting has been shifted to the function heatmap() which has now been defined as an S4 generic method. Previous methods for plotting heatmaps using plot() have been partly available in Versions 1.3.0 and 1.3.1. Since Version 1.3.2, they are no longer available.

## 11.2 Upgrading to Version 1.3.3 or newer

Users who upgrade to Version 1.3.3 (or newer) from an older version should be aware that the package now requires a newer version of Rcpp. This issue can simply be solved by re-installing Rcpp from CRAN using install.packages("Rcpp").

## 11.3 Upgrading to Version 1.4.0

The function sparseToFull() has been deprecated. A fully compatible function as.DenseSimilarityMatrix() is available that replaces and extends sparseToFull().

12 Change Log 57

# 12 Change Log

### Version 1.4.3:

 added optional color legend to heatmap plotting; in line with this change, some minor changes to the interface of the heatmap() function

• corresponding updates of help pages and vignette

#### **Version 1.4.2:**

- switched sequence kernel example in vignette from kernlab to kebabs package
- workaround to ensure that all apcluster\*() methods are able to process KernelMatrix objects (cf. kebabs package)
- replaced data set ch22Promoters by plain text file (FASTA format) in inst/examples
- bug fix in the heatmap() method
- vignette engine changed from Sweave to knitr

### **Version 1.4.1:**

■ fixes in C++ code of sparse affinity propagation

### Version 1.4.0:

- added apcluster() method for sparse similarity matrices; as a consequence, the package now imports the Matrix package and is now also able to handle non-sparse matrix classes defined by the Matrix package. Moreover, similarity functions supplied to the apcluster() method may now also return any matrix type defined by the Matrix package.
- fix of apcluster() for dense matrices to better support -Inf similarities
- added apclusterK() method for sparse similarity matrices
- preferenceRange() is now an S4 generic; re-implementation in C++ to speed up function; changed handling of -Inf similarities for consistency with sparse version.
- added preferenceRange() methods for sparse matrices and dense matrix objects from the Matrix package.
- new conversion methods implemented for converting dense similarity matrices to sparse ones and vice versa; consequently, sparseToFull() is marked as deprecated.
- adaptation of heatmap() function for improved handling of -Inf similarities
- adaptations of signatures of [ and [ [ accessor methods
- renamed help page of methods for computing similarity matrices to 'similarities' in order to avoid confusion with the accessor method similarity().
- corresponding updates of help pages and vignette

## **Version 1.3.5:**

memory access fixes in C++ code called from apclusterL()

58 12 Change Log

■ minor updates of vignette

### **Version 1.3.4:**

- added sort() function to rearrange clusters according to sort criterion; note that this
  is an S3 method (see help page for explanation)
- improvements and bug fixes of apclusterL() method for signature (s="matrix", x="missing")
- performance optimizations of apcluster() and apclusterL()
- plotting of clustering results superimposed in scatter plot matrices now also works for 'AggExResult' objects
- improvements of consistency of error and warning messages
- according adaptations of documentation and vignette
- adapted dependency and linking to Rcpp version 0.11.1 (to avoid issues on Mac OS)
- minor correction of package namespace

### Version 1.3.3:

- adapted dependencies and linking to Rcpp version 0.11.0
- cleared up package dependencies

### **Version 1.3.2:**

- plotting of clustering results extended to data sets with more than two dimensions (resulting in the clustering result being superimposed in a scatterplot matrix); the variant that plot() can be used to draw a heatmap has been removed. From now on, heatmap() must always be used.
- improved NA handling
- correction of input check in apcluster() and apclusterL() (previously, both functions issued a warning whenever argument p had length > 1)
- corresponding updates and further improvements of help pages and vignette

## **Version 1.3.1:**

- re-implementation of heatmap() method: dendrograms can now be plotted even for APResult and ExClust objects as well as for cluster hierarchies based on prior clusterings; color bars can now be switched off and colors can be changed by user (by new sideColor argument); dendrograms can be switched on and off (by Rowv and Colv arguments);
- added as.hclust() and as.dendrogram() methods
- added new arguments base, showSamples, and horizto the plot() method with signature (x="AggExResult", y="missing"); moreover, parameters for changing the appearance of the height axis are now respected as well
- streamlining of methods (redundant definition of inherited methods removed)
- various minor improvements of code and documentation

12 Change Log 59

### **Version 1.3.0:**

- added Leveraged Affinity Propagation Clustering
- re-implementation of main functions as S4 generic methods in order to facilitate the convenient internal computation of similarity matrices
- for convenience, similarity matrices can be stored as part of clustering results
- heatmap plotting now done by heatmap() which has been defined as S4 generic
- extended interface to functions for computing similarity matrices
- added function corSimMat()
- implementation of length() method for classes APResult, AggExResult, and ExClust
- added accessor function to extract clustering levels from AggExResult objects
- correction of exemplars returned by apcluster for details=TRUE in slot idxAll of returned APResult object
- when using data stored in a data frame, categorical columns are now explicitly omitted, thereby, avoiding warnings
- all clustering methods now store their calls into the result objects
- updates and extensions of help pages and vignette

### **Version 1.2.0:**

- reimplementation of apcluster() in C++ using the Rcpp package [3] which reduces computation times by a factor of 9-10
- obsolete function apclusterLM() removed
- updates of help pages and vignette

### **Version 1.1.1:**

- updated citation
- minor corrections in help pages and vignette

# Version 1.1.0:

- added exemplar-based agglomerative clustering function aggExCluster()
- added various plotting functions for dendrograms and heatmaps
- extended help pages and vignette according to new functionality
- added sequence analysis example to vignette along with data set ch22Promoters
- re-organization of variable names in vignette
- added option verbose to apclusterK()
- numerous minor corrections in help pages and vignette

## Version 1.0.3:

- Makefile in inst/doc eliminated to avoid installation problems
- renamed vignette to "apcluster.Rnw"

60 References

### **Version 1.0.2:**

replacement of computation of responsibilities and availabilities in function apcluster() by pure matrix operations (see 10.6 above); traditional implementation à la Frey and Dueck still available as function apclusterLM;

- improved support for named objects (see 10.2)
- new function for computing label vectors (see 10.3)
- re-organization of package source files and help pages

**Version 1.0.1:** first official release, released March 2, 2010

# 13 How to Cite This Package

If you use this package for research that is published later, you are kindly asked to cite it as follows:

U. Bodenhofer, A. Kothmeier, and S. Hochreiter (2011). APCluster: an R package for affinity propagation clustering. *Bioinformatics* **27**(17):2463–2464. DOI: 10.1093/bioinformatics/btr406.

Moreover, we insist that, any time you cite the package, you also cite the original paper in which affinity propagation has been introduced [5].

To obtain BibT<sub>E</sub>X entries of the two references, you can enter the following into your R session:

```
toBibtex(citation("apcluster"))
```

# References

- [1] U. Bodenhofer, A. Kothmeier, and S. Hochreiter. APCluster: an R package for affinity propagation clustering. *Bioinformatics*, 27(17):2463–2464, 2011.
- [2] B. De Baets and R. Mesiar. Metrics and *T*-equalities. *J. Math. Anal. Appl.*, 267:531–547, 2002.
- [3] D. Eddelbuettel and R. François. Rcpp: seamless R and C++ integration. *J. Stat. Softw.*, 40(8):1–18, 2011.
- [4] C. H. FitzGerald, C. A. Micchelli, and A. Pinkus. Functions that preserve families of positive semidefinite matrices. *Linear Alg. Appl.*, 221:83–102, 1995.
- [5] B. J. Frey and D. Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972–976, 2007.
- [6] A. K. Jain, M. N. Murty, and P. J. Flynn. Data clustering: a review. *ACM Comput. Surv.*, 31(3):264–323, 1999.

References 61

[7] A. Karatzoglou, A. Smola, K. Hornik, and A. Zeileis. kernlab – an S4 package for kernel methods in R. *J. Stat. Softw.*, 11(9):1–20, 2004.

- [8] C. Leslie, E. Eskin, and W. S. Noble. The spectrum kernel: a string kernel for SVM protein classification. In R. B. Altman, A. K. Dunker, L. Hunter, K. Lauderdale, and T. E. D. Klein, editors, *Pacific Symposium on Biocomputing* 2002, pages 566–575. World Scientific, 2002.
- [9] C. A. Micchelli. Interpolation of scattered data: Distance matrices and conditionally positive definite functions. *Constr. Approx.*, 2:11–22, 1986.
- [10] R. S. Michalski and R. E. Stepp. Clustering. In S. C. Shapiro, editor, *Encyclopedia of artificial intelligence*, pages 168–176. John Wiley & Sons, Chichester, 1992.
- [11] J. Palme, S. Hochreiter, and U. Bodenhofer. KeBABS: an R package for kernel-based analysis of biological sequences. *Bioinformatics*, 31(15):2574–2576, 2015.
- [12] B. Schölkopf and A. J. Smola. *Learning with Kernels*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2002.
- [13] J. H. Ward Jr. Hierarchical grouping to optimize an objective function. *J. Amer. Statist. Assoc.*, 58:236–244, 1963.