

The BClustLonG Package: A Dirichlet Process Mixture Model for Clustering Longitudinal Gene Expression Data

Jiehuan Sun¹, Jose D. Herazo-Maya², Naftali Kaminski², Hongyu Zhao¹, and Joshua L. Warren¹

¹*Department of Biostatistics, Yale School of Public Health*

²*Pulmonary, Critical Care and Sleep Medicine, Yale School of Medicine*

Overview

Many clustering methods have been proposed, but most of them cannot work for longitudinal gene expression data. Our newly developed method, BClustLonG, can be used to perform clustering analysis for longitudinal gene expression data. It adopts a linear-mixed effects framework to model the trajectory of genes over time, while clustering is jointly conducted based on the regression coefficients obtained from all genes. To account for the correlations among genes and alleviate the high dimensionality challenges, factor analysis models are adopted for the regression coefficients. The Dirichlet process prior distribution is utilized for the means of the regression coefficients to induce clustering (See Sun et al. (2017) for details).

This document provides a tutorial for using the **BClustLonG** package. The tutorial includes information on (1) the format of the input data and (2) how to obtain clustering results and visually show the clustering structure. As with any R package, detailed information on functions, along with their arguments and values, can be obtained in the help files.

Input data format

The analyses performed in this tutorial are based on a simulated dataset, which comes with the package. Basically, the data are generated from a mixture of two multivariate normal distributions, for which the covariance matrix satisfies the factor analysis model assumption (See simulation studies in Sun et al. (2017) for details).

The input data for **BClustLonG** has to be a list with three elements: Y (gene expression data), ID, and years (The names of the elements have to be matched exactly). Each column of Y represents one gene. The j_{th} row of “Y” represents the gene expression value for subject $ID[j]$ at time $years[j]$. So, the length of the “ID”, the length of the “years”, and the number of rows of “Y” are the same. No missing values are allowed in the data and the variable “years” is preferably on the scale of one (e.g. if the visiting time is on the scale of years, then the unit should be in year such as 2 years instead of 730 days or 24 months), which is related to the default values for the hyperparameters.

```
library(BClustLonG, quietly=TRUE)
data(data)
str(data)
```

```
## List of 3
## $ Y      : num [1:250, 1:10] 1.835 2.055 2.926 0.698 1.922 ...
##   ..- attr(*, "dimnames")=List of 2
##     .. ..$ : NULL
##     .. ..$ : chr [1:10] "Gene1" "Gene2" "Gene3" "Gene4" ...
## $ ID     : int [1:250] 1 1 1 1 1 2 2 2 2 2 ...
```


Obtain clustering results with desired number of clusters

Although *BClustLonG* can automatically determine the number of clusters and corresponding cluster structure, it is possible to utilize the outputs of *BClustLonG* to produce clustering structure for a given number of clusters (this feature could be useful in practice, since researchers might have some ideas on the number of clusters that are clinically meaningful). Specifically, provided the posterior similarity matrix, Hierarchical Clustering method (*hclust*) can be used to produce clustering structure for a given number of clusters. The following code shows how to achieve this, where the desired number of clusters is 4.

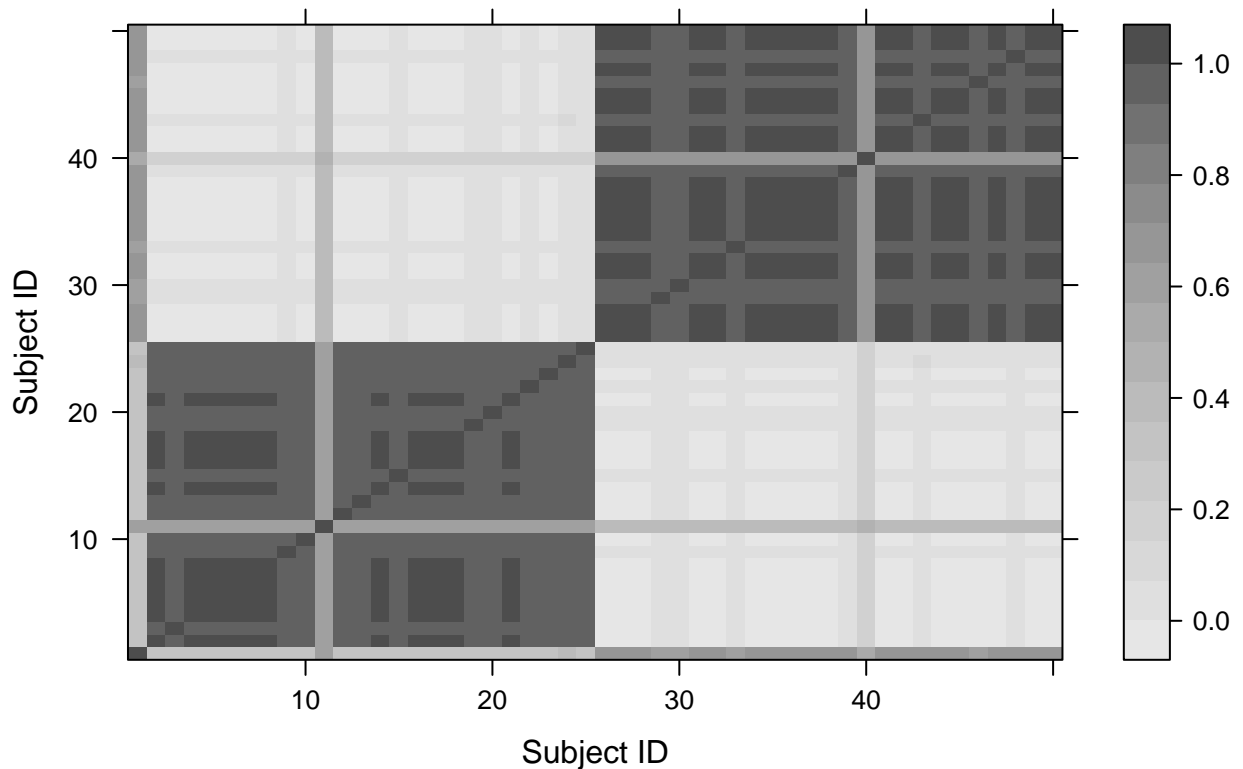
```
## using Hierarchical Clustering method to obtain the clustering results ##
CL = cutree(hclust(as.dist(1-mat)),k=4)
CL

## [1] 1 2 2 2 2 2 2 2 2 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 4 4 4 4 4 4 4 4 4 4
## [36] 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
```

Visualization of the similarity matrix

The posterior similarity matrix can also be visually shown so that users can have a rough idea of how many clusters there are. As shown in the figure, it is clear that there are two major clusters.

```
## plot similarity matrix ##
require(lattice,quietly=TRUE)
n = length(unique(data$ID))
x = rep(1:n,times=n)
y = rep(1:n,each=n)
z = as.vector(mat)
levelplot(z~x*y,col.regions=rev(gray.colors(n^2)), xlab = "Subject ID",ylab = "Subject ID")
```



Other options

BClustLonG adopts a linear-mixed effects with random intercepts and slopes to model the trajectory of genes over time. Hence, each subject has a subject-specific intercept and slope for each gene. It is possible that only the intercepts (i.e. baseline gene expression profiles) are informative of the clustering structure, in which case we may want to cluster only on the intercepts. This can be done by setting the parameter *infoVar*="int". Sometimes, a diagonal covariance matrix instead of the factor analysis model is preferred for the intercepts and slopes and this can be done by setting the parameter *factor*=TRUE (this can speed up the algorithm and may only be used to get an initial sense of the possible clustering structure in the data).

```
## Clustering based only on intercepts ##
res = BClustLonG(data, iter=500, thin=2,savePara=FALSE, infoVar="int",factor=TRUE)
## clustering based on intercepts and slopes ##
## assume diagonal covariance matrix for the intecepts and slopes ##
res = BClustLonG(data, iter=500, thin=2,savePara=FALSE, infoVar="both",factor=FALSE)
```

References

Sun, Jiehuan, Jose D Herazo-Maya, Naftali Kaminski, Hongyu Zhao, and Joshua L Warren. 2017. "A Dirichlet Process Mixture Model for Clustering Longitudinal Gene Expression Data." *Statistics in Medicine* 36 (22). Wiley Online Library: 3495–3506.