

Clusternomics

Stephen Coleman

17/10/2019

Introduction

Principal aim: to model both **global** (across dataset) and **local** (dataset-specific) clustering structure. Specifically, Clusternomics (Gabasova, Reid, and Wernisch 2017) can allow for clusters merging and separating in different datasets (also referred to as **contexts**). A motivating example is shown below for two 1D datasets.

```
library(ggplot2) # ubiquitous
library(ggExtra) # for ggMarginal

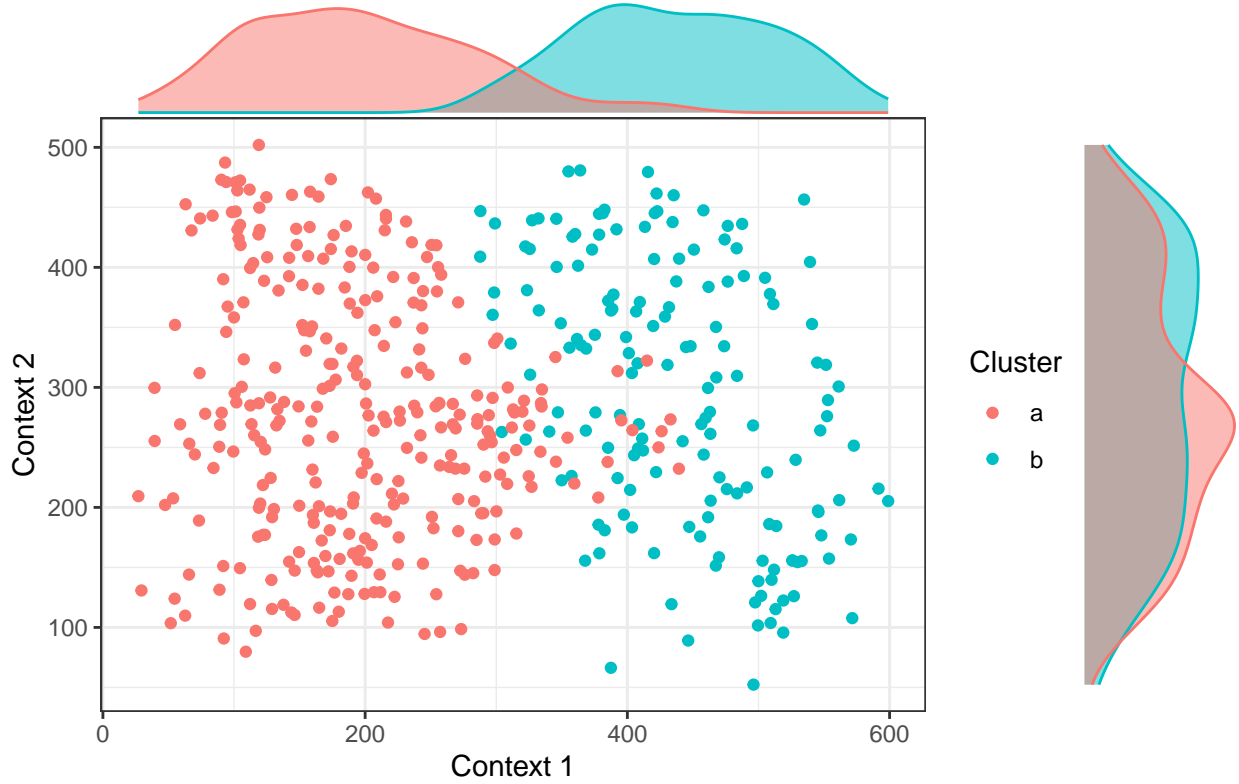
# Personal preference
theme_set(theme_bw())

# Data generated at http://drawdata.xyz/
my_data <- read.csv("./example_data.csv")

# Plot the data
p1 <- ggplot(data = my_data, mapping = aes(x = x, y = y, colour = z)) +
  geom_point() +
  labs(
    title = "Example of different cluster behaviour across contexts",
    x = "Context 1",
    y = "Context 2",
    colour = "Cluster"
  )

# Add marginal density plots by grouping
ggMarginal(p1, groupColour = T, groupFill = T)
```

Example of different cluster behaviour across contexts



Here we have two separate subpopulations. In Context 1 the sub-populations are discernible as two local clusters. However this not the case in Context 2. Clusternomics allows for such disagreement in local models while conveying the complexity of clustering structure (in this case that there really is two sub-populations present) to the global model.

The model

Clusternomics, which is embedded in the Bayesian clustering framework, uses a hierarchical Dirichlet mixture model to identify structure on both the local and the global level. The model is fit to the data via Gibbs sampling. The model does not assume that cluster behaviour will be consistent across heterogeneous datasets. This is not to assume that the clustering structure uncovered in one dataset should not inform the clustering in another dataset. This can be summarised as so:

1. Clustering structure in one dataset should inform the clustering in another dataset. If two points are clustered together in one context they should be more inclined to cluster together in other contexts.
2. Different degrees of dependence should be allowed between clusters across contexts. The model should work when datasets have the same underlying structure and also when each dataset is effectively independent of all others. Fundamental to this is allowing datasets to have different numbers of clusters.

To enable these modelling aims, Clusternomics explicitly represents the local clusters (i.e. dataset specific) and the global structure that emerges when considering the combination of the datasets. The global clusters are defined by combinations of local clusters. Consider the case where 3 clusters emerge in Context 1 (denoted by labels $\{1, 2, 3\}$) and 2 clusters emerge in Context 2 (denoted by labels $\{A, B\}$). In this case our global structure has the possible form:

$$\{(1, A), (2, A), (3, A), (1, B), (2, B), (3, B)\}$$

Thus if a point is assigned a label of 1 in Context 1 and a label of A in Context 2 it increases the probability of cluster $(1, A)$ becoming populated at the global level. However, it is possible that some of the possible global clusters described above are not realised as some local clusters overlap across datasets. Consider the case that labellings 1 and 2 from the first context are captured entirely by label A in the second context with a label of 3 corresponding perfectly to label B . In this case our global structure would take the form:

$$\{(1, A), (2, A), (3, B)\}$$

In this way the local structure informs the global structure.

The original paper introduces two models that are “asymptotically equivalent”. The first is easier to develop an intuition of, but it is the second that is implemented as it is more computationally efficient.

Notation

Let us denote the number of datasets by L and all observed data by X . Let

$$\begin{aligned} X &= (X_1, \dots, X_L), \\ X_l &= (X_{l1}, \dots, X_{ln}) \end{aligned}$$

where X_l is the data of the l th context.

It is assumed that it is the same n individuals in each dataset in the same order. Therefore we have L membership vectors denoting cluster membership:

$$\begin{aligned} C &= (C_1, \dots, C_L), \\ C_l &= (c_{l1}, \dots, c_{ln}). \end{aligned}$$

Basic model

The basis of the integrative model is a finite approximation of a Dirichlet process known as a Dirichlet-Multinomial Allocation mixture model (Green and Richardson 2001). There is a nice explanation of this model in Savage et al. (2013).

In this case we model the latent structure in the l th dataset using a mixture of K_l components. This means that the full model density is the weighted sum of the probability density functions associated with each component where the weights, π_{lk} , are the proportion of the total population assigned to the k th components:

$$\begin{aligned} p(X_{li}|c_{li} = k) &= \pi_{lk} f(X_{li}|\theta_{lk}), \\ p(X_{li}) &= \sum_{k=1}^{K_l} \pi_{lk} f_l(X_{li}|\theta_{lk}). \end{aligned}$$

Here K_l is the total number of clusters present and θ_{lk} are the parameters defining the k th distribution in the l th dataset.

The weights, $\pi_l = (\pi_{l1}, \dots, \pi_{lK_l})$, follow a Dirichlet distribution with concentration parameter α_0 .

The distributions in the mixture model for each dataset are:

$$\begin{aligned}
\pi_l &\sim \text{Dirichlet}\left(\frac{\alpha_0}{K_l}, \dots, \frac{\alpha_0}{K_l}\right) \\
c_{li} &\sim \text{Categorical}(\pi_l) \\
\theta_{lk} &\sim h_l \\
X_{li}|c_{li} = k &\sim f_l(X_{li}|\theta_{lk})
\end{aligned}$$

where H_l is some prior distribution for parameters for each mixture component; F_l is a probability distribution for samples given the parameters θ_{lk} . Note that each context may have different distributions and hyperparameters.

First formulation: the intuitive model

For ease of understanding, let $L = 2$. Then each context has its own mixture weights with symmetric Dirichlet priors:

$$\begin{aligned}
\pi_1 &\sim \text{Dirichlet}\left(\frac{\alpha_1}{K_1}, \dots, \frac{\alpha_1}{K_1}\right) \\
\pi_2 &\sim \text{Dirichlet}\left(\frac{\alpha_2}{K_2}, \dots, \frac{\alpha_2}{K_2}\right)
\end{aligned}$$

These weights form the basis of the local clustering within each dataset. A third mixture distribution is used to link the two local clusters. This has a Dirichlet prior over the global mixture weights, ρ , which is defined over the outer product of the local weights:

$$\rho \sim \text{Dirichlet}(\gamma \text{vec}(\pi_1 \otimes \pi_2)).$$

The outer product is defined:

$$\begin{aligned}
\pi_1 \otimes \pi_2 &= \pi_1 \pi_2^T \\
&= \begin{pmatrix} \pi_{11}\pi_{21} & \pi_{11}\pi_{22} & \cdots & \pi_{11}\pi_{2K_2} \\ \pi_{12}\pi_{21} & \pi_{12}\pi_{22} & \cdots & \pi_{12}\pi_{2K_2} \\ \vdots & \vdots & \ddots & \vdots \\ \pi_{1K_1}\pi_{21} & \pi_{1K_1}\pi_{22} & \cdots & \pi_{1K_1}\pi_{2K_2} \end{pmatrix}
\end{aligned}$$

Then the vector function takes the column vectors of the matrix and places them in one vector:

$$\text{vec}(\pi_1 \otimes \pi_2) = \begin{pmatrix} \pi_{11}\pi_{21} \\ \vdots \\ \pi_{1K_1}\pi_{21} \\ \pi_{11}\pi_{22} \\ \vdots \\ \pi_{1K_1}\pi_{22} \\ \pi_{11}\pi_{23} \\ \vdots \\ \pi_{1K_1}\pi_{2K_2} \end{pmatrix}$$

This is the basis for the non-symmetric concentration parameter of the Dirichlet distribution over the global mixture weights, ρ .

A global membership variable, c , is drawn from a Categorical distribution with concentration parameter ρ :

$$\begin{aligned} c_i &= (c_{1i}, c_{2i}) \\ c &\sim \text{Categorical}(\rho) \\ X_{li}|c_{li} = k &\sim f_l(X_{li}|\theta_{lk}), l \in \{1, 2\} \end{aligned}$$

In this way the local clusters, c_{li} , in the l th context are projections of the global clustering, c_i , onto this space.

The model achieves the property that the local assignment should affect the global assignment posterior probability. This can be seen as local assignments affect the posterior probability of the context-specific weights, π_l , which in turn define the probabilities of the global combinatorial clusters through the hierarchical model (a change in the membership of the k th component in the l th dataset affects an entire slice of column vectors in the tensors defining the concentration parameter for the global component weights, ρ).

The model also represents different degrees of dependence by allowing any combination of cluster assignments across contexts. When there is a single common cluster structure across the two contexts, the occupied clusters will be concentrated along the diagonal of the probability matrix of ρ .

When $L > 2$ then the outer product of the local component weights generalises to a tensor product. In this case the hierarchical model is given by:

$$\begin{aligned} \pi_l|\alpha_0 &\sim \text{Dirichlet}\left(\frac{\alpha_0}{K_l}, \dots, \frac{\alpha_0}{K_l}\right) \\ \rho|\gamma, \{\pi_1 \dots, \pi_L\} &\sim \text{Dirichlet}\left(\gamma \left\{ \bigotimes_{l=1}^L \pi_l \right\}\right) \\ c_i|\rho &\sim \text{Categorical}(\rho), c_i = (c_{1i}, \dots, c_{Li}) \\ \theta_{lk} &\sim h_l \\ X_{li}|c_{li} = k &\sim f_l(X_{li}|\theta_{lk}) \end{aligned}$$

The model for L contexts has a scaling issue as each additional dataset brings an additional layer of calculations in the probability tensor. The second model is designed to reduce this cost by avoiding calculations for uninhabited components.

Second formulation: the quick model

This model attempts to reduce the number of combinations required. This requires decoupling the number of local clusters and the number of global clusters. A mixture over S global clusters is defined:

$$\begin{aligned} \rho|\gamma_0 &\sim \text{Dirichlet}\left(\frac{\gamma_0}{S}, \dots, \frac{\gamma_0}{S}\right) \\ c_i|\rho &\sim \text{Categorical}(\rho) \end{aligned}$$

where ρ is the global mixture weights and c_i is the components assignment indicator variable as before. A variable z_{ls} is then defined that associates the s th global cluster with context specific clusters:

$$\begin{aligned} \pi_l|\alpha_0 &\sim \text{Dirichlet}\left(\frac{\alpha_0}{K_l}, \dots, \frac{\alpha_0}{K_l}\right), \\ z_{ls}|\pi_l &\sim \text{Categorical}(\pi_l). \end{aligned}$$

Here $z_{ls} \in \{1, \dots, K_l\}$ assigns the s th global cluster to the a specific local cluster in the l th dataset. One may consider $(z_{sl})_{s=1}^S$ as the coordinates of the global clusters in the l th dataset. This link means that the c_i maps a point to a global cluster as well as the local clusters.

In the previous model the mapping of global clusters to local clusters was implicit, because each combination of context clusters mapped to a unique global cluster. In this model, the mapping is probabilistic and forms a part of the model. One may now state S , the number of global components, and thus limit the number of computations required in contrast to the preceding model.

This hierarchial model is defined by:

$$\begin{aligned}
\rho|\gamma_0 &\sim \text{Dirichlet}\left(\frac{\gamma_0}{S}, \dots, \frac{\gamma_0}{S}\right), \\
c_i|\rho &\sim \text{Categorical}(\rho), \\
\pi_l|\alpha_0 &\sim \text{Dirichlet}\left(\frac{\alpha_0}{K_l}, \dots, \frac{\alpha_0}{K_l}\right), \\
z_{ls}|\pi_l &\sim \text{Categorical}(\pi_l) \\
\theta_{kl}|\pi_l &\sim h_l \\
X_{li}|c_i, (z_{sl})_{s=1}^S, (\theta_{lk})_{k=1}^K &\sim f_l(X_{li}|\theta_{lk_{c_i}})
\end{aligned}$$

Inference

The quick, decoupled model is implemented in the R package `clusternomics` and uses Gibbs sampling as the inference algorithm.

R package

Setup

The `clusternomics` package is available on CRAN.

```

# If not installed, remedy this.
if (!requireNamespace("clusternomics", quietly = TRUE)) {
  devtools::install_github("evelinag/clusternomics")
  # install.packages("clusternomics")
}

# Load the library
library(clusternomics)

# For the pipe
library(magrittr)
library(plyr)

```

Evelina has a vignette showcasing the Clusternomics model on simulated data at <https://cran.r-project.org/web/packages/clusternomics/vignettes/using-clusternomics.html>.

First, define two local clusters in each of two contexts:

$$\begin{aligned}
\text{Cluster 1: } x &\sim \mathcal{N}\left(\begin{pmatrix} -1.5 \\ -1.5 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) \\
\text{Cluster 2: } x &\sim \mathcal{N}\left(\begin{pmatrix} 1.5 \\ 1.5 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right)
\end{aligned}$$

Here the first row of each vector is the parameters for the normal distribution in the first context and the second row is the value of the parameters in the second context.

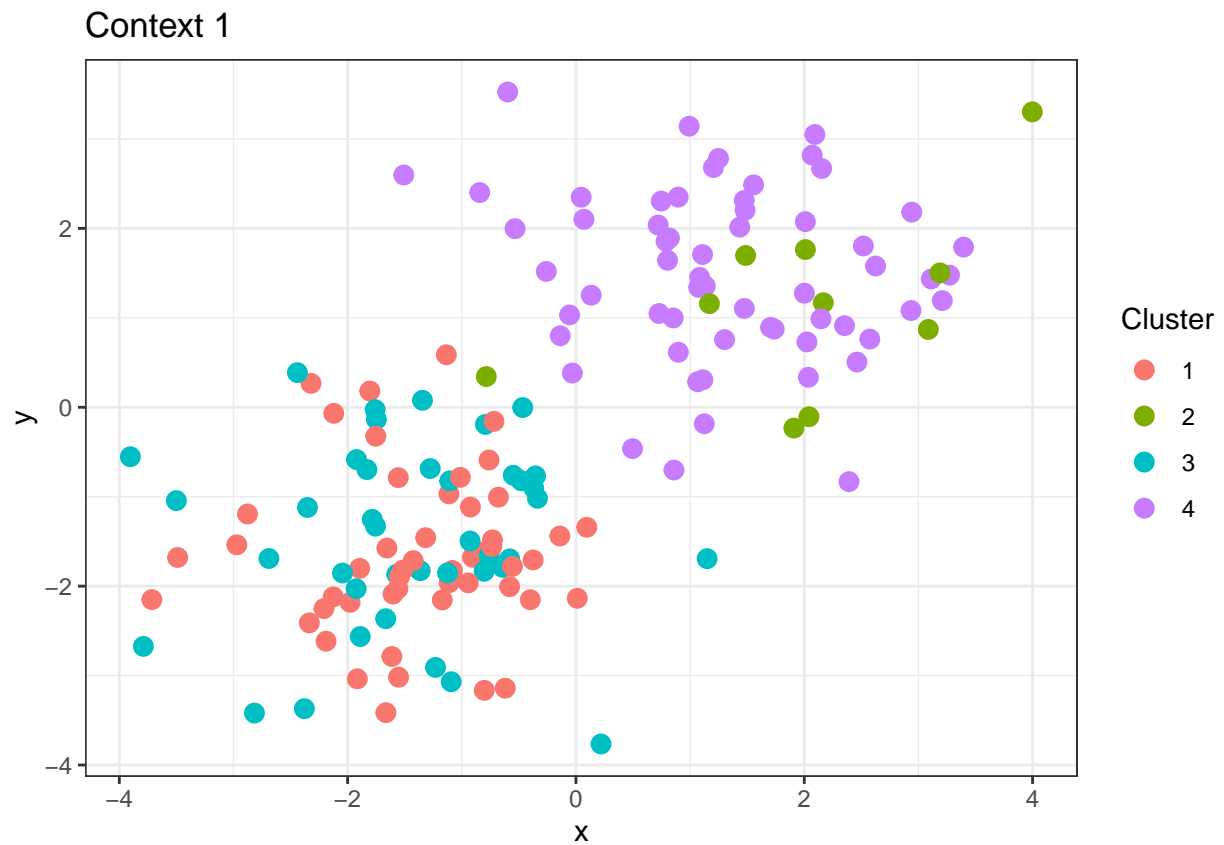
	Cluster 1 (Context 2)	Cluster 2 (Context 2)
Cluster 1 (Context 1)	50	10
Cluster 2 (Context 1)	40	60

There are 160 data points in total. There are two local clusters within each context. At the same time, there are four clusters on the global level that appear when we combine observations from the local contexts.

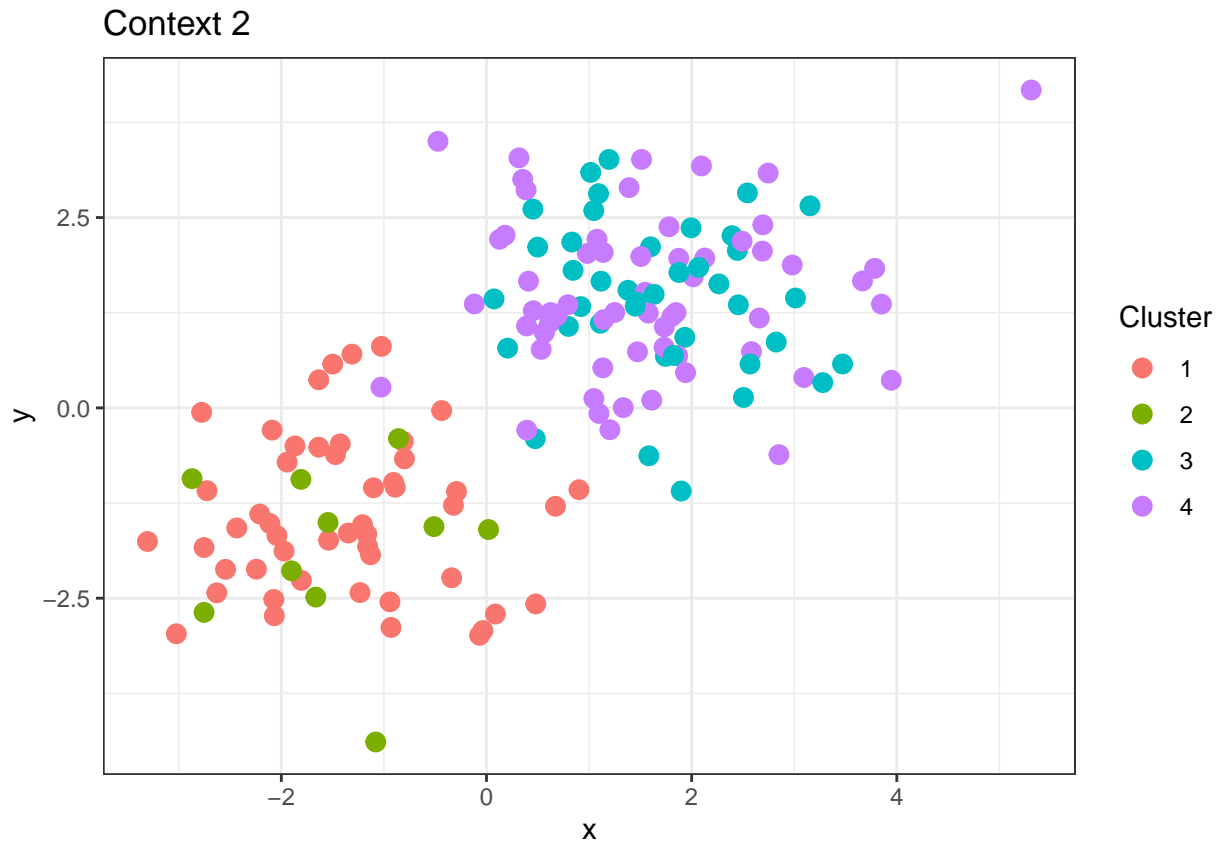
```
# For consistent results
set.seed(1)

# Number of elements in each cluster, follows the table given above
groupCounts <- c(50, 10, 40, 60)
# Centers of clusters
means <- c(-1.5, 1.5)
# Helper function to generate test data
testData <- generateTestData_2D(groupCounts, means)
datasets <- testData$data

# Inspect the data in each context
# Context 1
qplot(datasets[[1]][, 1], datasets[[1]][, 2], col = factor(testData$groups)) +
  geom_point(size = 3) +
  ggtitle("Context 1") + xlab("x") + ylab("y") +
  scale_color_discrete(name = "Cluster")
```



```
# Context 2  
qplot(datasets[[2]][, 1], datasets[[2]][, 2], col = factor(testData$groups)) +  
  geom_point(size = 3) +  
  ggtitle("Context 2") + xlab("x") + ylab("y") +  
  scale_color_discrete(name = "Cluster")
```

Note that the following assumptions work awfully well in this case as the data is simulated. In real life the number of MCMC iterations would have to be significantly larger to see convergence-like behaviour. We can also use relatively low upper bounds on the number of clusters as we know the ground truth. This also decreases the time required within each individual iteration.

```
# Setup of the algorithm
dataDistributions <- "diagNormal"

# Pre-specify number of clusters
clusterCounts <- list(global = 10, context = c(3, 3))

# Set number of iterations
# The following is ONLY FOR SIMULATION PURPOSES
# Use larger number of iterations for real-life data
maxIter <- 300
burnin <- 200
lag <- 2 # Thinning of samples

# Run context-dependent clustering
results <- contextCluster(datasets, clusterCounts,
  maxIter = maxIter,
  burnin = burnin,
  lag = lag,
  dataDistributions = "diagNormal",
  verbose = F
)

# Extract resulting cluster assignments
```

```

samples <- results$samples

# Extract global cluster assignments for each MCMC sample
clusters <- lapply(1:length(samples), function(i) samples[[i]]$Global)
# clusters <- sapply(1:length(samples), function(i) samples[[i]]$Global) %>%
# t()

```

Model analysis

Convergence

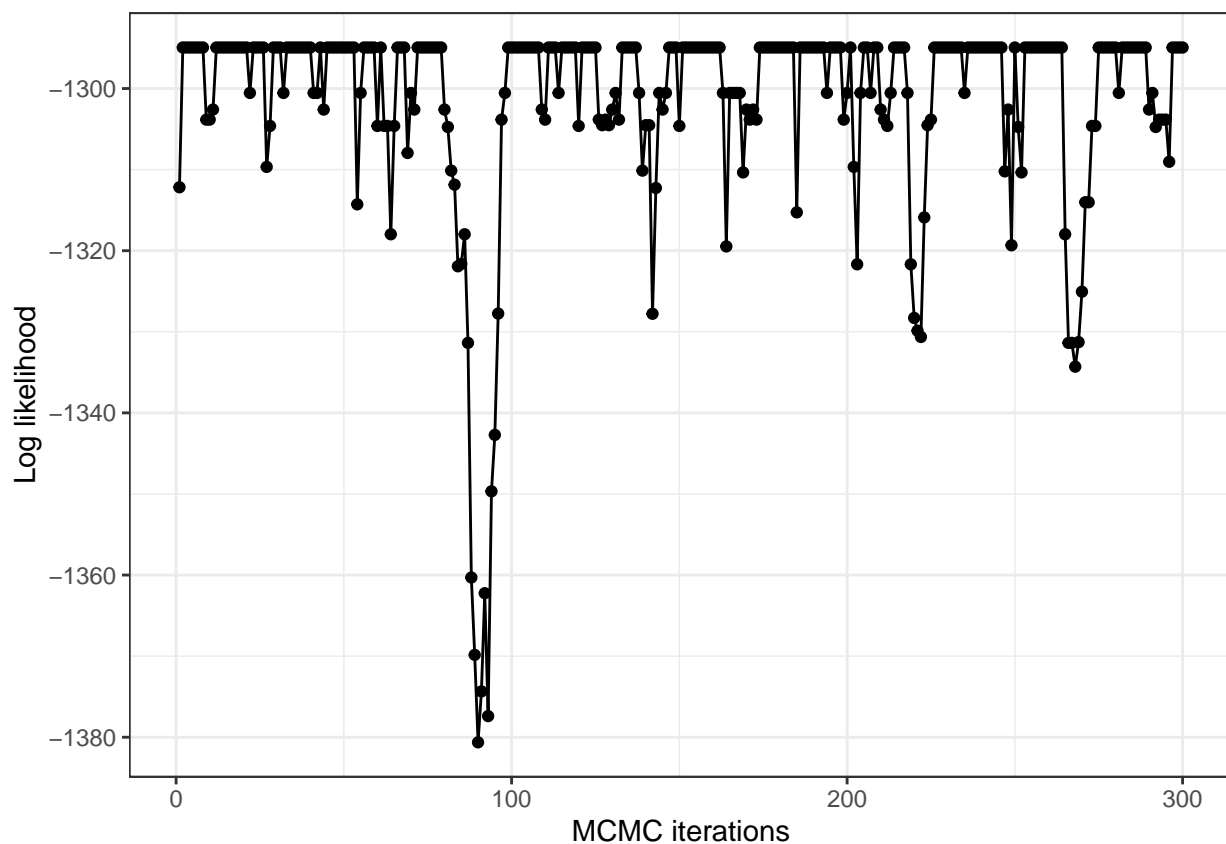
The default for checking convergence, trace plots. In this case, Gabasova, Reid, and Wernisch (2017) recommend log-likelihood as the quantity to measure.

```
logliks <- results$logliks
```

```

qplot(1:maxIter, logliks) +
  geom_line() +
  xlab("MCMC iterations") +
  ylab("Log likelihood")

```



The below is copy and pasted from <https://cran.r-project.org/web/packages/clusternomics/vignettes/using-clusternomics.html> as I don't think there's much to say here.

Choosing number of clusters

As part of the training, we also compute the Deviance Information Criterion (DIC). The DIC is a Bayesian model selection method that can be used to select the number of clusters. The DIC value is returned in `results$DIC`. Models that better fit the data will result in lower values of DIC than worse models. For example, if we fit a model with number of clusters that is too small, we get higher DIC value than for the original result.

```
wrongClusterCounts <- list(global = 2, context = c(2, 1))
worseResults <- contextCluster(datasets, wrongClusterCounts,
  maxIter = maxIter,
  burnin = burnin,
  lag = lag,
  dataDistributions = "diagNormal",
  verbose = F
)

print(paste("Original model has lower (better) DIC:", results$DIC))

## [1] "Original model has lower (better) DIC: 2604.21806172573"

print(paste("Worse model has higher (worse) DIC:", worseResults$DIC))

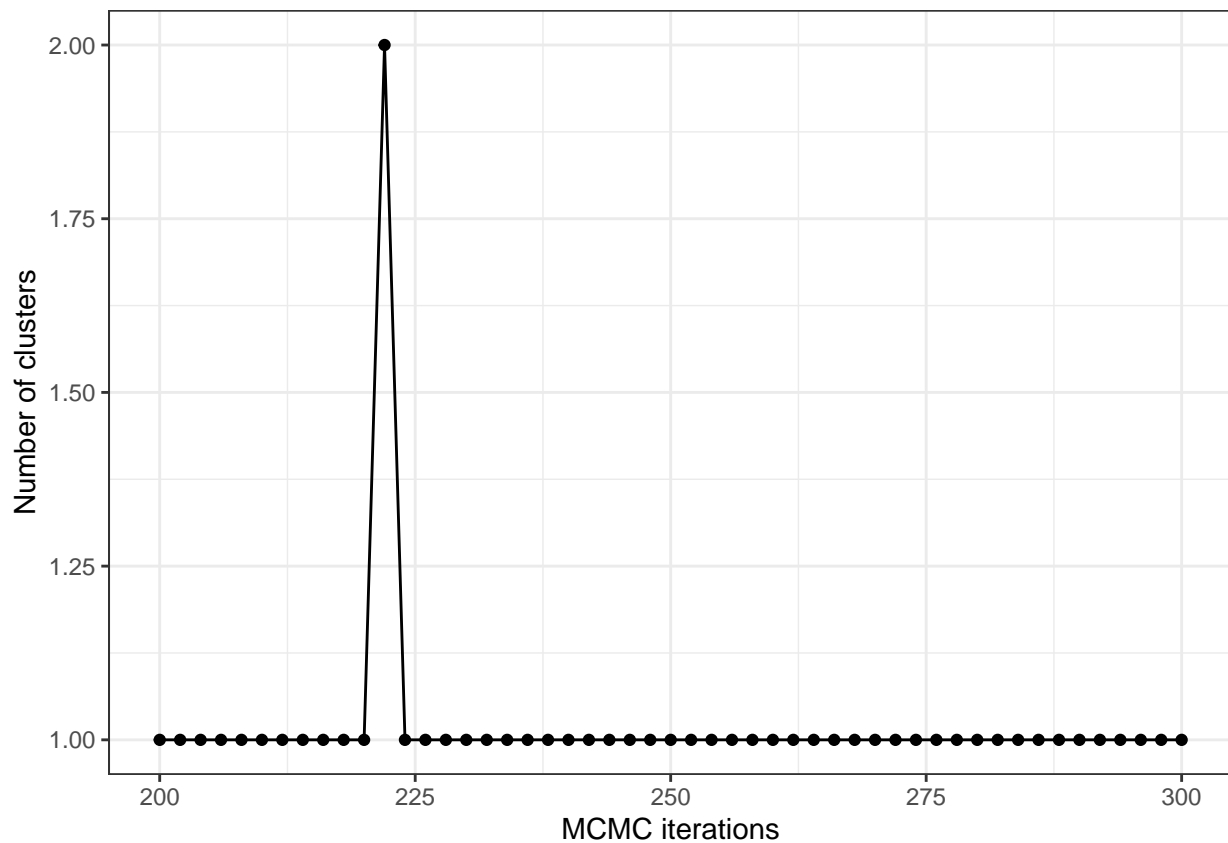
## [1] "Worse model has higher (worse) DIC: 2584.25139051911"
```

Posterior number of clusters

We can also look at the number of global clusters that were identified in the datasets. The plot below shows the number of global clusters across MCMC samples. This is the number of actually occupied global clusters, which can be smaller than the number of global clusters specified when running the `contextCluster` function.

```
# This does not run as the number of clusters found is precisely 1.
cc <- numberOfClusters(clusters) # This function does not seem to work
# cc <- lapply(1:length(clusters), function(i) unique(clusters[[i]]) %>% length )

qplot(seq(from = burnin, to = maxIter, by = lag), cc) +
  geom_line() +
  xlab("MCMC iterations") +
  ylab("Number of clusters")
```



Sizes of clusters

Here we look at the posterior sizes of the individual global clusters across the MCMC iterations and then we show a box plot with the estimated sizes. The labels of global clusters represent the corresponding combinations of local clusters.

*# There's some issues in the code - I suspect that certain outputs where changed
and the input of other functions never updated.*

This issue was dplyr::laply. Now I use that.

Find all the cluster labels in the retained MCMC samples

```
clusterLabels <- unique(clusters %>% as.vector())
```

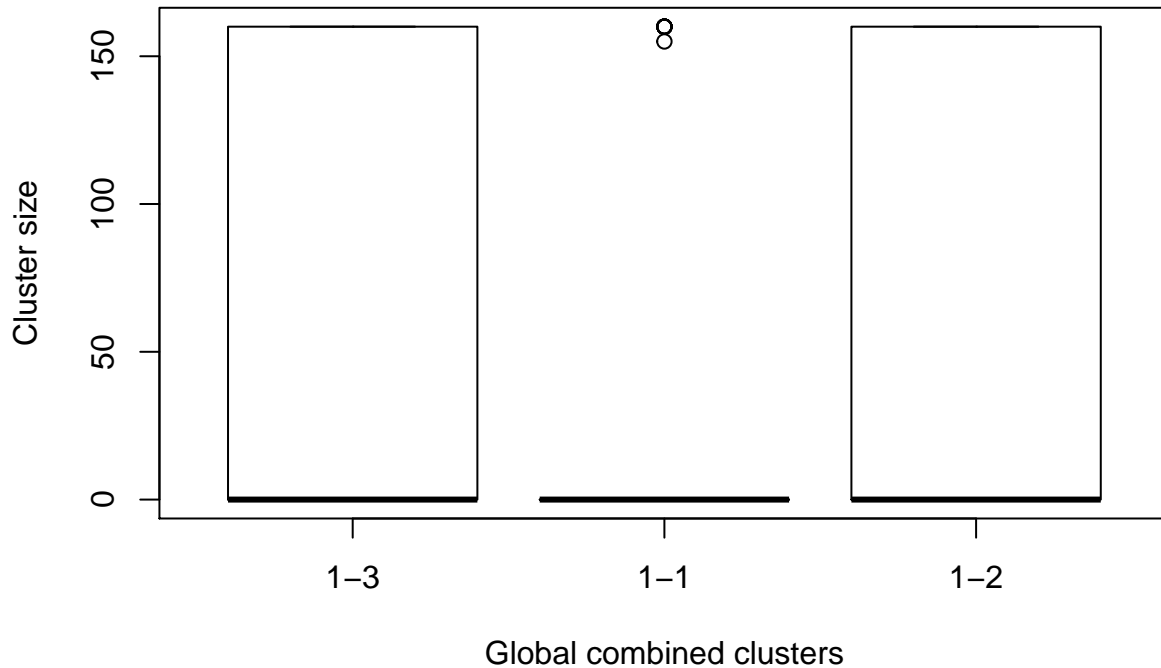
```
sizes <- matrix(nrow = nrow(clusters), ncol = length(clusterLabels))
```

```
for (ci in 1:length(clusterLabels)) {
  sizes[, ci] <- rowSums(clusters == clusterLabels[ci])
}
```

```
sizes <- sizes %>% as.data.frame()
```

```
colnames(sizes) <- clusterLabels
```

```
boxplot(sizes, xlab = "Global combined clusters", ylab = "Cluster size")
```



```
# # Create a matrix to store the cluster membership count in each sample
# sizes <- matrix(nrow = length(clusters), ncol = length(clusterLabels))
#
# for(ii in 1:length(clusters)){
#   for (ci in 1:length(clusterLabels)) {
#     sizes[ii, ci] <- sum(clusters[[ii]] == clusterLabels[[ci]])
#   }
# }
# sizes <- sizes %>% as.data.frame()
# colnames(sizes) <- clusterLabels
#
# # Plot
# boxplot(sizes, xlab = "Global combined clusters", ylab = "Cluster size")
```

Obtaining global clusters

There are several approaches to estimate hard cluster assignments from MCMC samples. If the MCMC chain converged to a stable results, it is possible to use one of the samples as the resulting cluster assignment.

```
clusteringResult <- samples[[length(samples)]]
```

Co-clustering matrix

A more principled approach is to look at which data points were assigned into the same cluster across the MCMC samples. We can explore this using the posterior co-clustering matrix, which estimates the posterior probability that two samples belong to the same cluster.

```
# Compute the co-clustering matrix from global cluster assignments
coclust <- coclusteringMatrix(clusters)
```

```

# Plot the co-clustering matrix as a heatmap
require(gplots)

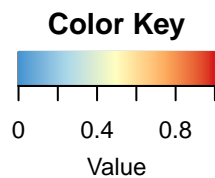
## Loading required package: gplots

##
## Attaching package: 'gplots'

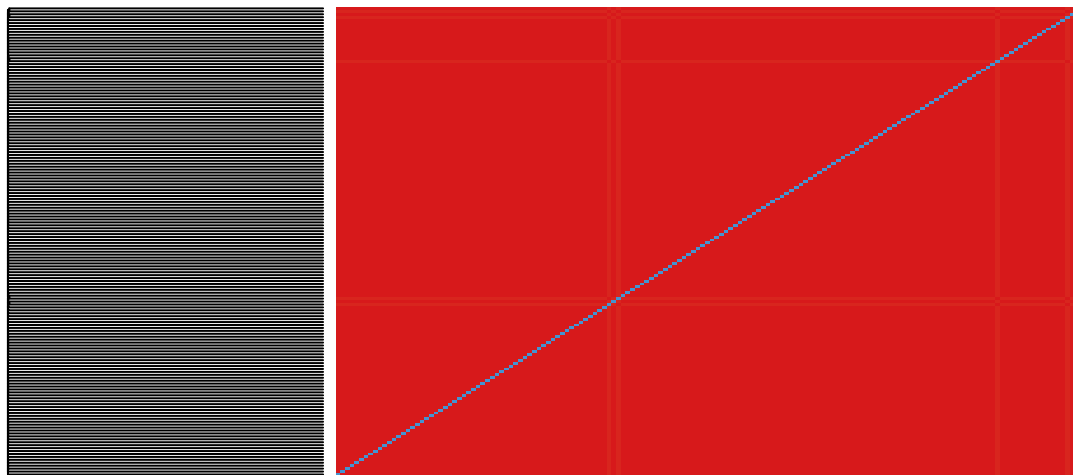
## The following object is masked from 'package:stats':
##
##      lowess

mypalette <- colorRampPalette(rev(c("#d7191c", "#fdae61", "#ffffbf", "#abd9e9", "#4395d2")),
  space = "Lab"
)(100)
h <- heatmap.2(
  coclust,
  col = mypalette, trace = "none",
  dendrogram = "row", labRow = "", labCol = "", key = TRUE,
  keysize = 1.5, density.info = c("none"),
  main = "MCMC co-clustering matrix",
  scale = "none"
)

```



MCMC co-clustering matrix



We can then use the posterior co-clustering matrix to compute the resulting hard clustering using hierarchical clustering. Note that for this step, we need to specify the number of clusters that we want to obtain. This step should be guided by the posterior number of clusters estimated by the algorithm (see above).

```

diag(coclust) <- 1
fit <- hclust(as.dist(1 - coclust))
hardAssignments <- cutree(fit, k = 4)

```

Adjusted Rand index

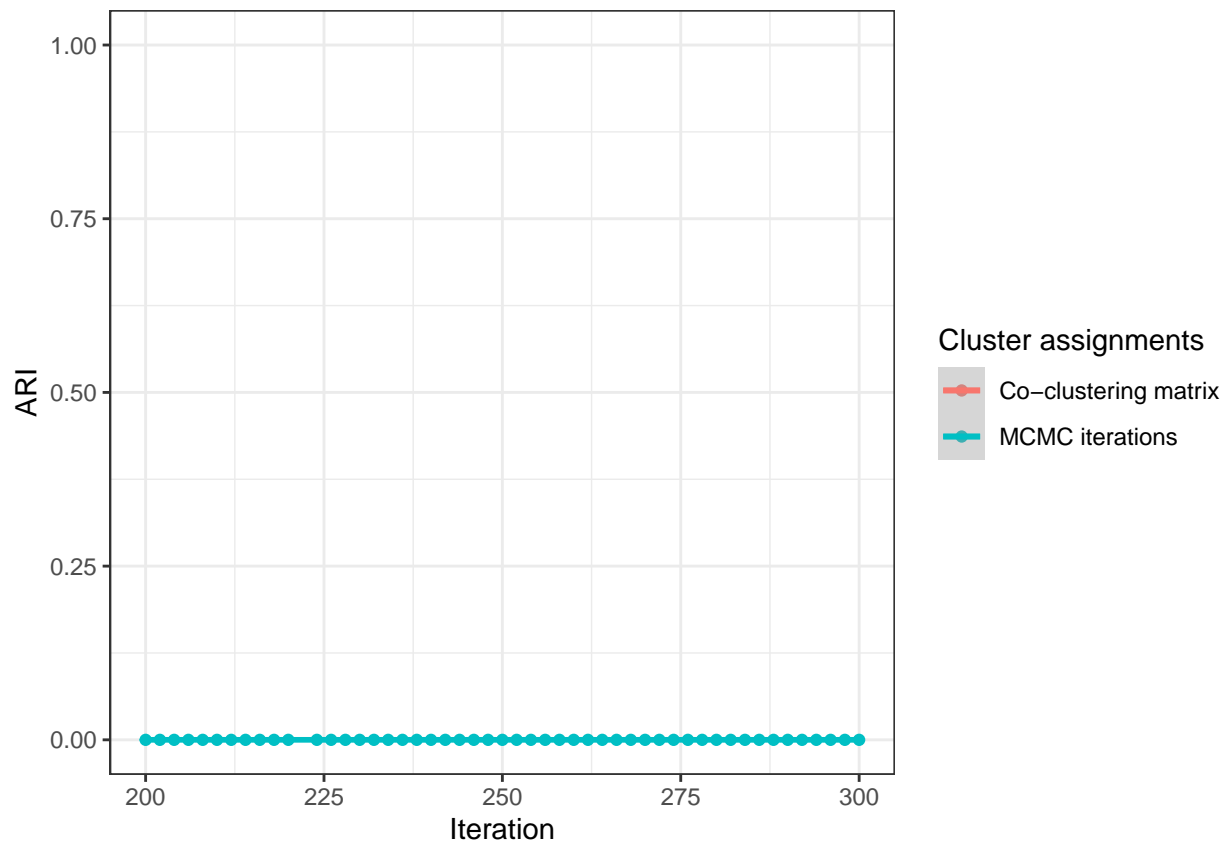
Now we can check if the estimated global cluster assignments correspond to the true assignments that were used to generate the simulated dataset. We use the adjusted Rand index (ARI), which measures how well do two sets of assignments correspond to each other. The following plot shows the ARI across the MCMC iterations, and also the ARI of the result obtained from the co-clustering matrix. ARI equal to 1 represents complete agreement between the estimated assignments and the true clustering, ARI equal to 0 corresponds to random assignments.

```
# Calculate the adjusted rand index
aris <- laply(
  1:nrow(clusters),
  function(i) {
    mclust::adjustedRandIndex(clusters[i, ], testData$groups)
  }
) %>%
  as.data.frame()

colnames(aris) <- "ARI"
aris$Iteration <- seq(from = burnin, to = maxIter, by = lag)
coclustAri <- mclust::adjustedRandIndex(hardAssignments, testData$groups)
aris$Coclust <- coclustAri

ggplot(aris, aes(x = Iteration, y = ARI, colour = "MCMC iterations")) +
  geom_point() +
  ylim(0, 1) +
  geom_smooth(size = 1) +
  theme_bw() +
  geom_line(aes(x = Iteration, y = Coclust, colour = "Co-clustering matrix"), size = 1) +
  scale_colour_discrete(name = "Cluster assignments")

## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
## Warning: Removed 1 rows containing non-finite values (stat_smooth).
## Warning: Removed 1 rows containing missing values (geom_point).
## Warning: Removed 51 rows containing missing values (geom_path).
```



References

- Gabasova, Evelina, John Reid, and Lorenz Wernisch. 2017. "Clusternomics: Integrative Context-Dependent Clustering for Heterogeneous Datasets." *PLoS Computational Biology* 13 (10). Public Library of Science: e1005781.
- Green, Peter J, and Sylvia Richardson. 2001. "Modelling Heterogeneity with and Without the Dirichlet Process." *Scandinavian Journal of Statistics* 28 (2). Wiley Online Library: 355–75.
- Savage, Richard S, Zoubin Ghahramani, Jim E Griffin, Paul Kirk, and David L Wild. 2013. "Identifying Cancer Subtypes in Glioblastoma by Combining Genomic, Transcriptomic and Epigenomic Data." *arXiv Preprint arXiv:1304.3577*.