

Consensus clustering for Bayesian mixture models: Supplementary materials

Stephen Coleman, Paul DW Kirk and Chris Wallace

September 11, 2020

Abstract

- 1 Algorithms**
- 2 The model**
- 3 Additional results**
 - 3.1 Simulations**
 - 3.2 Yeast**

Input: Distance between means Δ_μ
A common standard deviation σ^2
A number of clusters K
The number of items to generate in total N
The number of features to generate in total P
An indicator vector of feature relevance $\phi = (\phi_1, \dots, \phi_P)$
The expected proportion of items in each cluster $\pi = (\pi_1, \dots, \pi_K)$
A method for sampling x times from the array y , with weights π : *Sample*(y, x, π)
A method for permuting a vector x : *Permute*(x)
A method for generating a value from a univariate Gaussian distribution with mean μ and standard deviation σ^2 : *Gaussian*(μ, σ^2)
Output: A dataset, X
The generating cluster labels $c = (c_1, \dots, c_N)$

```

begin
  /* initialise the empty data matrix */
   $X \leftarrow 0_{N \times P}$ ;
  /* create a matrix of  $K$  means */
   $\mu \leftarrow (\Delta_\mu, \dots, K\Delta_\mu)$ ;
  /* generate the allocation vector */
   $c \leftarrow \text{Sample}(1 : K, N, \pi)$ ;
   $M \leftarrow 0_{N \times N}$ ;
  for  $p = 1$  to  $P$  do
    /* Test if the feature is relevant, if relevant
       generate data from a mixture of univariate
       Gaussians, otherwise draw all items from the same
       distribution */
    if  $\phi_p = 1$  then
       $\nu \leftarrow \text{Permute}(\mu)$ ;
      for  $n = 1$  to  $N$  do
         $X(n, p) \leftarrow \text{Gaussian}(\nu_{c_n}, \sigma^2)$ 
      end
    end
    if  $\phi_p = 0$  then
      for  $n = 1$  to  $N$  do
         $X(n, p) \leftarrow \text{Gaussian}(0, \sigma^2)$ 
      end
    end
  end
  /* Mean centre and scale the data */
   $X \leftarrow \text{Normalise}(X)$ 
end

```

Algorithm 1: Data generation for a mixture of Gaussian with independent features. This algorithm is implemented in the `generateSimulationDataset` function from the `mdiHelpR` package available at www.github.com/stcolema/mdiHelpR.

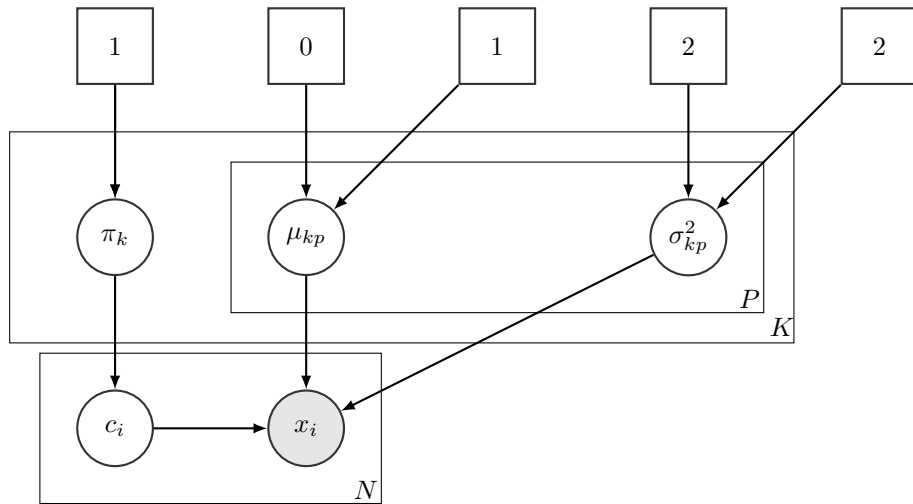


Figure 1: Directed acyclic graph for the Bayesian mixture of Gaussians used.

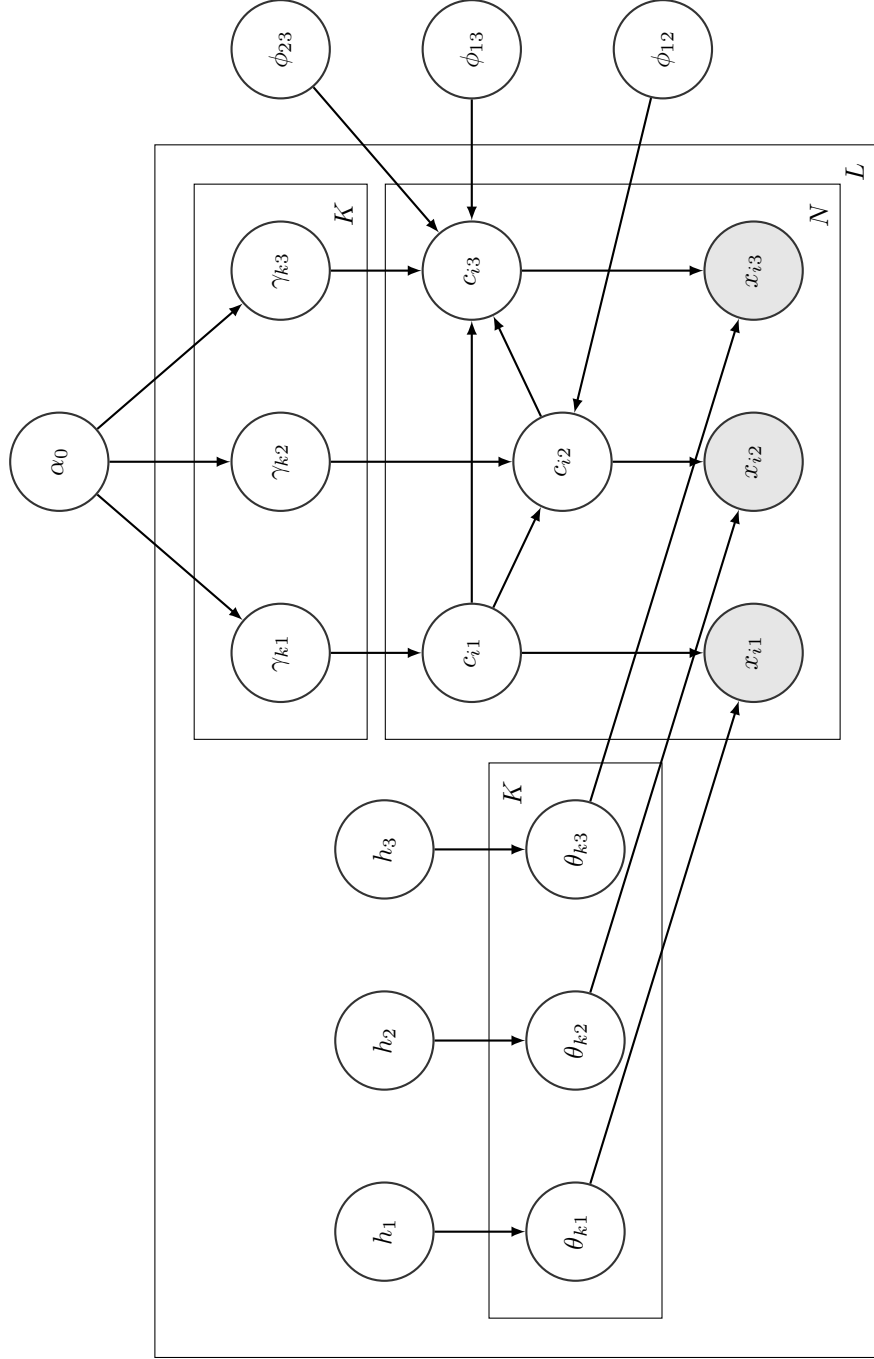


Figure 2: Directed acyclic graph for the Multiple Dataset Integration model for 3 datasets.

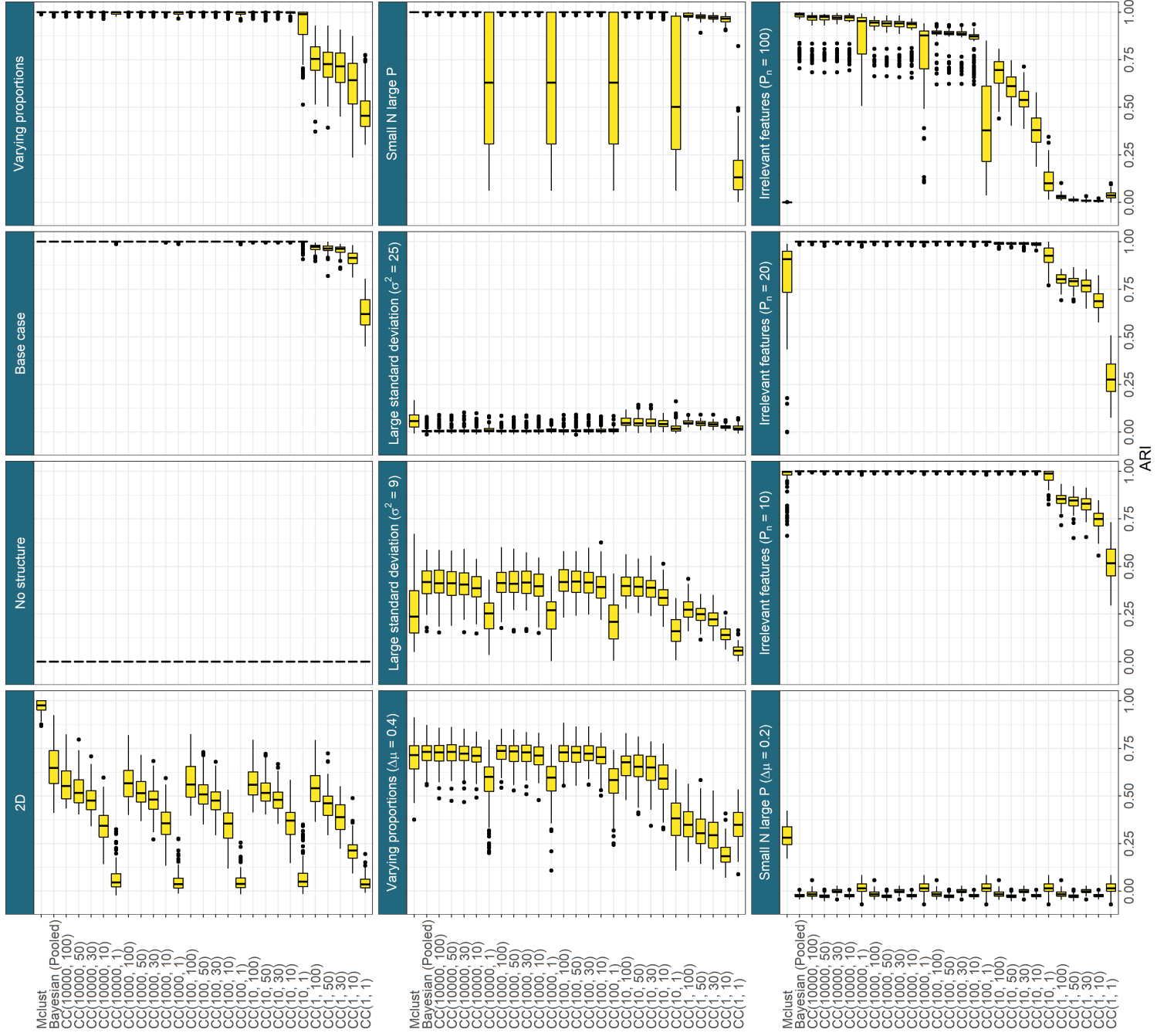


Figure 3: Predictive performance across all simulations. $CC(R, S)$ denotes consensus clustering using the R^{th} sample from S different chains.

Uncertainty quantification

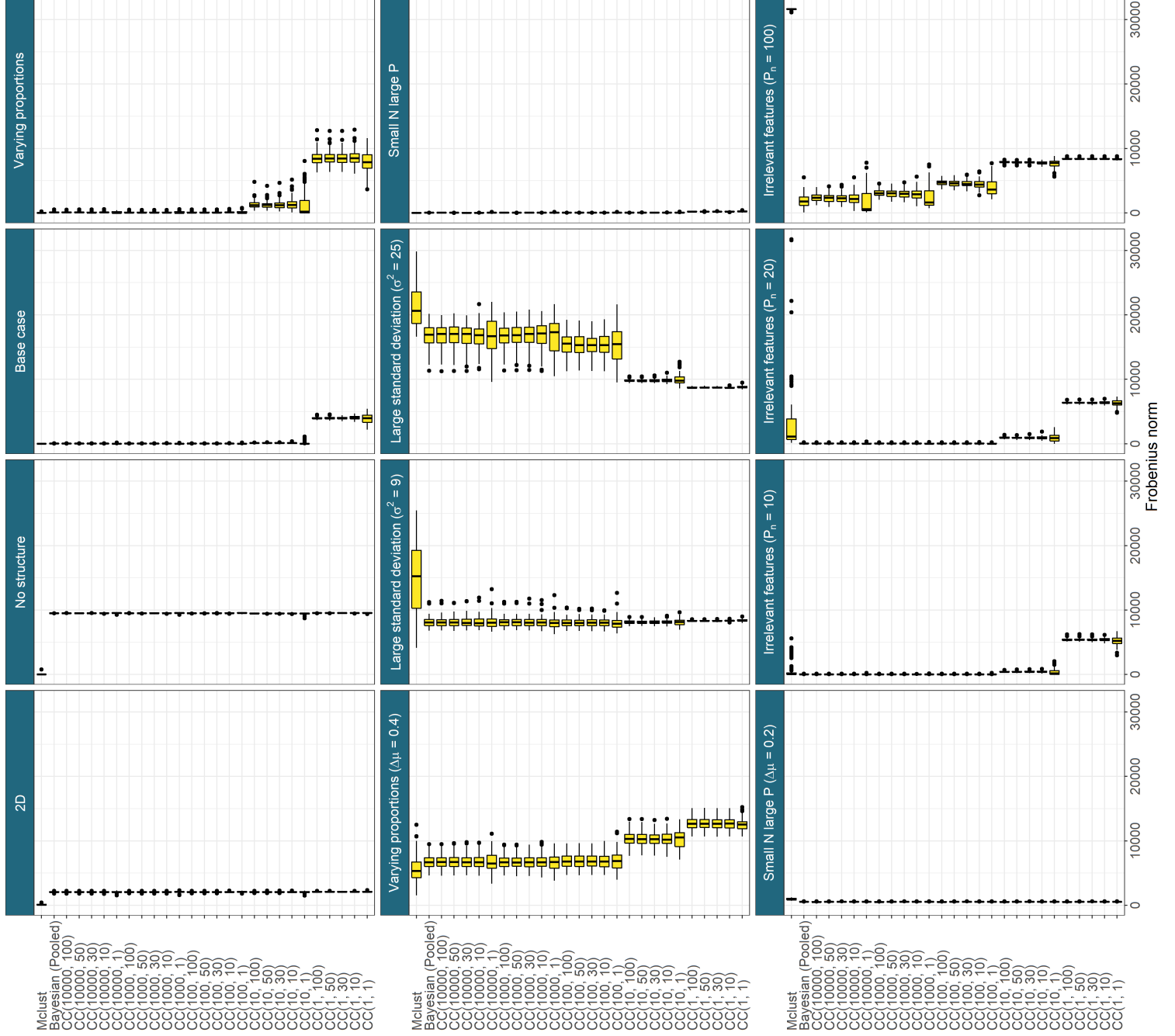


Figure 4: Frobenius norm across simulations. $CC(R, S)$ denotes consensus clustering using the R^{th} sample from S different chains. Lower values are better.



Figure 5: The GO enrichment is very similar across chains.