
CONSENSUS CLUSTERING FOR BAYESIAN MIXTURE MODELS

Stephen Coleman^{1,*}

stephen.coleman@mrc-bsu.cam.ac.uk

Paul D.W. Kirk^{1,2,†}

paul.kirk@mrc-bsu.cam.ac.uk

Chris Wallace^{1,2,†}

cew54@cam.ac.uk

¹ MRC Biostatistics Unit

² Cambridge Institute of Therapeutic Immunology & Infectious Disease

University of Cambridge, U.K.

* Corresponding author.

† These authors provided an equal contribution.

ABSTRACT

1 Cluster analysis is an integral part of precision medicine and systems biology, used to define
2 groups of patients or biomolecules. Consensus clustering is an ensemble approach that
3 is widely used in these areas, which combines the output from multiple runs of a non-
4 deterministic clustering algorithm. Here we consider the application of consensus clustering
5 to a broad class of heuristic clustering algorithms that can be derived from Bayesian mixture
6 models (and extensions thereof) by adopting an early stopping criterion when performing
7 sampling-based inference for these models. While the resulting approach is non-Bayesian, it

Consensus clustering for Bayesian mixture models

8 inherits the usual benefits of consensus clustering, particularly in terms of computational
9 scalability and providing assessments of clustering stability/robustness.

10 In simulation studies, we show that our approach can successfully uncover the target clus-
11 tering structure, while also exploring different plausible clusterings of the data. We show
12 that, when a parallel computation environment is available, our approach offers significant
13 reductions in runtime compared to performing sampling-based Bayesian inference for the
14 underlying model, while retaining many of the practical benefits of the Bayesian approach,
15 such as exploring different numbers of clusters. We propose a heuristic to decide upon
16 ensemble size and the early stopping criterion, and then apply consensus clustering to a
17 clustering algorithm derived from a Bayesian integrative clustering method. We use the
18 resulting approach to perform an integrative analysis of three 'omics datasets for budding
19 yeast and find clusters of co-expressed genes with shared regulatory proteins. We validate
20 these clusters using data external to the analysis. These clusters can help assign likely
21 function to understudied genes, for example *GAS3* clusters with histones active in S-phase,
22 suggesting a role in DNA replication.

23 Our approach can be used as a wrapper for essentially any existing sampling-based Bayesian
24 clustering implementation, and enables meaningful clustering analyses to be performed using
25 such implementations, even when computational Bayesian inference is not feasible, e.g. due
26 to poor scalability. This enables researchers to straightforwardly extend the applicability of
27 existing software to much larger datasets, including implementations of sophisticated models
28 such as those that jointly model multiple datasets.

29 **Keywords** Cluster analysis · Multiomics · Ensemble learning

30 **Background**

31 From defining a taxonomy of disease to creating molecular sets, grouping items can help us to understand
32 and make decisions using complex biological data. For example, grouping patients based upon disease
33 characteristics and personal omics data may allow the identification of more homogeneous subgroups,
34 enabling stratified medicine approaches. Defining and studying molecular sets can improve our understanding
35 of biological systems as these sets are more interpretable than their constituent members (1), and study of
36 their interactions and perturbations may have ramifications for diagnosis and drug targets (2, 3). The act
37 of identifying such groups is referred to as *cluster analysis*. Many traditional methods such as *K*-means

Consensus clustering for Bayesian mixture models

38 clustering (4, 5) condition upon a fixed choice of K , the number of clusters. These methods are often heuristic
39 in nature, relying on rules of thumb to decide upon a final value for K . For example, different choices of K
40 are compared under some metric such as silhouette (6) or the within-cluster sum of squared errors (**SSE**) as a
41 function of K . Moreover, K -means clustering can exhibit sensitivity to initialisation, necessitating multiple
42 runs in practice (7).

43 Another common problem is that traditional methods offer no measure of the stability or robustness of the
44 final clustering. Returning to the stratified medicine example of clustering patients, there might be individuals
45 that do not clearly belong to any one particular cluster; however if only a point estimate is obtained, this
46 information is not available. Ensemble methods address this problem, as well as reducing sensitivity to
47 initialisation. These approaches have had great success in supervised learning, most famously in the form of
48 Random Forest (8) and boosting (9). In clustering, consensus clustering (10) is a popular method which has
49 been implemented in R (11) and to a variety of methods (12, 13) and been applied to problems such as cancer
50 subtyping (14, 15) and identifying subclones in single cell analysis (16). Consensus clustering uses W runs
51 of some base clustering algorithm (such as K -means). These W proposed partitions are commonly compiled
52 into a *consensus matrix*, the $(i, j)^{th}$ entries of which contain the proportion of model runs for which the i^{th}
53 and j^{th} individuals co-cluster (for this and other definitions see section 1 of the Supplementary Material),
54 although this step is not fundamental to consensus clustering and there is a large body of literature aimed at
55 interpreting a collection of partitions (see, e.g., 17, 18, 19). This consensus matrix provides an assessment
56 of the stability of the clustering. Furthermore, ensembles can offer reductions in computational runtime
57 because the individual members of the ensemble are often computationally inexpensive to fit (e.g, because
58 they are fitted using only a subset of the available data) and because the learners in most ensemble methods
59 are independent of each other and thus enable use of a parallel environment for each of the quicker model
60 runs (20).

61 Traditional clustering methods usually condition upon a fixed choice of K , the number of clusters with the
62 choice of K being a difficult problem in itself. In consensus clustering, Monti *et al.* (10) proposed methods
63 for choosing K using the consensus matrix and Ünlü *et al.* (21) offer an approach to estimating K given
64 the collection of partitions, but each clustering run uses the same, fixed, number of clusters. An alternative
65 clustering approach, mixture modelling, embeds the cluster analysis within a formal, statistical framework
66 (22). This means that models can be compared formally, and problems such as the choice of K can be
67 addressed as a model selection problem (23). Moreover, *Bayesian mixture models* can be used to try to
68 directly infer K from the data. Such inference can be performed through use of a Dirichlet Process mixture
69 model (24, 25, 26), a mixture of finite mixture models (27, 28) or an over-fitted mixture model (29). These

Consensus clustering for Bayesian mixture models

70 models and their extensions have a history of successful application to a diverse range of biological problems
71 such as finding clusters of gene expression profiles (30), cell types in flow cytometry (31, 32) or scRNAseq
72 experiments (33), and estimating protein localisation (34). Bayesian mixture models can be extended to
73 jointly model the clustering across multiple datasets (35, 36) (section 2 of the Supplementary Material).

74 Markov chain Monte Carlo (MCMC) methods are the most common tool for performing computational
75 Bayesian inference. In Bayesian clustering, they are used to draw a collection of clustering partitions from
76 the posterior distribution. However, in practice, chains can become stuck in local posterior modes preventing
77 convergence (see, e.g., the Supplementary Materials of 37) and/or can require prohibitively long runtimes,
78 particularly when analysing high-dimensional datasets. Some MCMC methods make efforts to overcome
79 the problem of exploration, often at the cost of increased computational cost per iteration (38). There are
80 MCMC methods that use parallel chains to improve the scalability or reduce the bias of the Monte Carlo
81 estimate. However, these methods have various limitations. For instance, divide-and-conquer strategies
82 such as Asymptotically Exact, Embarrassingly Parallel MCMC (39) use subsamples of the dataset with each
83 chain to improve scaling with the number of items being clustered. This assumes that each subsample is
84 representative of the population, and is less helpful in situations where we have high-dimension but only
85 moderate sample size, such as analysis of 'omics data. Alternative approaches, such as distributed MCMC
86 (40) and coupling (41) have to account for burn-in bias; moreover, coupling further assumes the chains meet
87 in finite time and then stay together. In practice, a further challenge associated with these methods is that
88 their implementation may necessitate a substantial redevelopment of existing software.

89 Motivated by the lack of scalability of existing implementations of sampling-based Bayesian clustering (due
90 to prohibitive computational runtimes, as well as poor exploration, as described above), here we aim to
91 develop a general and straightforward procedure that exploits the flexibility of these methods, but extends
92 their applicability. Specifically, we make use of existing sampling-based Bayesian clustering implementations,
93 but only run them for a fixed (and relatively small) number of iterations, stopping before they have converged
94 to their target stationary distribution. Doing this repeatedly, we obtain an ensemble of clustering partitions,
95 which we use to perform consensus clustering. We propose a heuristic for deciding upon the ensemble size
96 (the number of learners used, W) and the ensemble depth (the number of iterations, D), inspired by the use
97 of scree plots in Principal Component Analysis (PCA; 42).

98 We show via simulation that our approach can successfully identify meaningful clustering structures. We then
99 illustrate the use of our approach to extend the applicability of existing Bayesian clustering implementations,
100 using as a case study the Multiple Dataset Iteration (MDI; 35) model for Bayesian integrative clustering
101 applied to real data. While the simulation results serve to validate our method, it is important to also evaluate

Consensus clustering for Bayesian mixture models

102 methods on real data which may represent more challenging problems. For our real data, we use three 'omics
103 datasets relating to the cell cycle of *Saccharomyces cerevisiae* with the aim of inferring clusters of genes
104 across datasets. As there is no ground truth available, we then validate these clusters using knowledge external
105 to the analysis.

106 Material and methods

107 Consensus clustering for Bayesian mixture models

108 We apply consensus clustering to MCMC based Bayesian clustering models using the method described in
algorithm 1. Our application of consensus clustering has two main parameters at the ensemble level, the

Data: $X = (x_1, \dots, x_N)$

Input:

The number of chains to run, W

The number of iterations within each chain, D

A clustering method that uses MCMC methods to generate samples of clusterings of the data

$Cluster(X, d)$

Output:

A predicted clustering, \hat{Y}

The consensus matrix M

begin

```
    /* initialise an empty consensus matrix */
```

```
     $M \leftarrow \mathbf{0}_{N \times N}$ ;
```

```
    for  $w = 1$  to  $W$  do
```

```
        /* set the random seed controlling initialisation and MCMC moves */
```

```
         $set.seed(w)$ ;
```

```
        /* initialise a random partition on  $X$  drawn from the prior distribution
```

```
           */
```

```
         $Y_{(0,w)} \leftarrow Initialise(X)$ ;
```

```
        for  $d = 1$  to  $D$  do
```

```
            /* generate a markov chain for the membership vector */
```

```
             $Y_{(d,w)} \leftarrow Cluster(X, d)$ ;
```

```
        end
```

```
        /* create a coclustering matrix from the  $D^{th}$  sample */
```

```
         $B^{(w)} \leftarrow Y_{(D,w)}$ ;
```

```
         $M \leftarrow M + B^{(w)}$ ;
```

```
    end
```

```
     $M \leftarrow \frac{1}{W} M$ ;
```

```
     $\hat{Y} \leftarrow$  partition  $X$  based upon  $M$ ;
```

```
end
```

Algorithm 1: Consensus clustering for Bayesian mixture models.

109

110 chain depth, D , and ensemble width, W . We infer a point clustering from the consensus matrix using the
111 maxpear function (43) from the R package mclust (44) which maximises the posterior expected adjusted

Consensus clustering for Bayesian mixture models

112 Rand index between the true clustering and point estimate if the matrix is composed of samples drawn from
113 the posterior distribution (section 3 of the Supplementary Material for details). There are alternative choices
114 of methods to infer a point estimate which minimise different loss functions (see, e.g., 45, 46, 47).

115 **Determining the ensemble depth and width**

116 As our ensemble sidesteps the problem of convergence within each chain, we need an alternative stopping rule
117 for growing the ensemble in chain depth, D , and number of chains, W . We propose a heuristic based upon
118 the consensus matrix to decide if a given value of D and W are sufficient. We suspect that increasing W and
119 D might continuously improve the performance of the ensemble, but we observe in our simulations that these
120 changes will become smaller and smaller for greater values, eventually converging for each of W and D . We
121 notice that this behaviour is analogous to PCA in that where for consensus clustering some improvement
122 might always be expected for increasing chain depth or ensemble width, more variance will be captured by
123 increasing the number of components used in PCA. However, increasing this number beyond some threshold
124 has diminishing returns, diagnosed in PCA by a scree plot. Following from this, we recommend, for some
125 set of ensemble parameters, $D' = \{d_1, \dots, d_I\}$ and $W' = \{w_1, \dots, w_J\}$, find the mean absolute difference
126 of the consensus matrix for the d_i^{th} iteration from w_j chains to that for the $d_{(i-1)}^{th}$ iteration from w_j chains
127 and plot these values as a function of chain depth, and the analogue for sequential consensus matrices for
128 increasing ensemble width and constant depth.

129 If this heuristic is used, we believe that the consensus matrix and the resulting inference should be stable
130 (see, e.g., 48, 49), providing a robust estimate of the clustering. In contrast, if there is still strong variation
131 in the consensus matrix for varying chain length or number, then we believe that the inferred clustering is
132 influenced significantly by the random initialisation and that the inferred partition is unlikely to be stable for
133 similar datasets or reproducible for a random choice of seeds.

134 **Simulation study**

We use a finite mixture with independent features as the data generating model within the simulation study. Within this model there exist “irrelevant features” (50) that have global parameters rather than cluster specific parameters. The generating model is

$$p(X, c, \theta, \pi | K) = p(K)p(\pi | K)p(\theta | K) \prod_{i=1}^N p(c_i | \pi, K) \prod_{p=1}^P p(x_{ip} | c_i, \theta_{c_i p})^{\phi_p} p(x_{ip} | \theta_p)^{(1-\phi_p)} \quad (1)$$

Consensus clustering for Bayesian mixture models

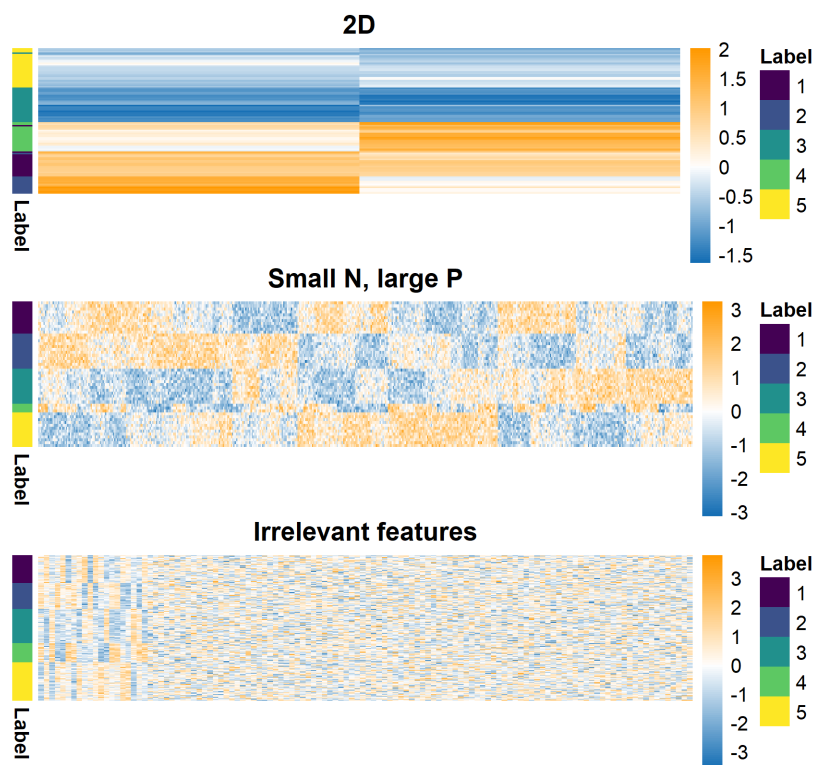


Figure 1: Example of generated datasets. Each row is an item being clustered and each column a feature of generated data. The 2D dataset (which is ordered by hierarchical clustering here) should enable proper mixing of chains in the MCMC. The small N , large P case has clear structure (observable by eye). This is intended to highlight the problems of poor mixing due to high dimensions even when the generating labels are quite identifiable. In the irrelevant features case, the structure is clear in the relevant features (on the left-hand side of this heatmap). This setting is intended to test how sensitive each approach is to noise.

135 for data $X = (x_1, \dots, x_N)$, cluster label or allocation variable $c = (c_1, \dots, c_N)$, cluster weight $\pi =$
 136 (π_1, \dots, π_K) , K clusters and the relevance variable, $\phi \in \{0, 1\}$ with $\phi_p = 1$ indicating that the p^{th} feature
 137 is relevant to the clustering. We used a *Gaussian* density, so $\theta_{kp} = (\mu_{kp}, \sigma_{kp}^2)$. We defined three scenarios
 138 and simulated 100 datasets in each (figure 1 and table 1). Additional details of the simulation process and
 139 additional scenarios are included in section 4.1 of the Supplementary Materials.

Table 1: Parameters defining the simulation scenarios as used in generating data and labels. $\Delta\mu$ is the distance between neighbouring cluster means within a single feature. The number of relevant features (P_s) is $\sum_p \phi_p$, and $P_n = P - P_s$.

Scenario	N	P_s	P_n	K	$\Delta\mu$	σ^2	π
2D	100	2	0	5	3.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Small N, large P	50	500	0	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Irrelevant features	200	20	100	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$

Consensus clustering for Bayesian mixture models

140 In each of these scenarios we apply a variety of methods (listed below) and compare the inferred point
141 clusterings to the generating labels using the Adjusted Rand Index (**ARI**, 51).

- 142 • `McLust`, a maximum likelihood implementation of a finite mixture of Gaussian densities (for a range
143 of modelled clusters, K),
- 144 • 10 chains of 1 million iterations, thinning to every thousandth sample for the overfitted Bayesian
145 mixture of Gaussian densities, and
- 146 • A variety of consensus clustering ensembles defined by inputs of W chains and D iterations within
147 each chain (see algorithm 1) with $W \in \{1, 10, 30, 50, 100\}$ and $D \in \{1, 10, 100, 1000, 10000\}$
148 where the base learner is an overfitted Bayesian mixture of Gaussian densities.

149 Note that none of the applied methods include a model selection step and as such there is no modelling of the
150 relevant variables. This and the unknown value of K is what separates the models used and the generating
151 model described in equation 1. More specifically, the likelihood of a point X_n for each method is

$$p(X_n|\mu, \Sigma, \pi) = \sum_{k=1}^K \pi_k p(X_n|\mu_k, \Sigma_k), \quad (2)$$

152 where $p(X_n|\mu_k, \Sigma_k)$ is the probability density function of the multivariate Gaussian distribution parameterised
153 by a mean vector, μ_k , and a covariance matrix, Σ_k , and π_k is the component weight such that $\sum_{k=1}^K \pi_k = 1$.
154 The implementation of the Bayesian mixture model restricts Σ_k to be a diagonal matrix while `McLust` models
155 a number of different covariance structures. Note that while we use the overfitted Bayesian mixture model,
156 this is purely from convenience and we expect that a true Dirichlet Process mixture or a mixture of mixture
157 models would display similar behaviour in an ensemble.

158 The ARI is a measure of similarity between two partitions, c_1, c_2 , corrected for chance, with 0 indicating c_1 is
159 no more similar to c_2 than a random partition would be expected to be and a value of 1 showing that c_1 and
160 c_2 perfectly align. Details of the methods in the simulation study can be found in sections 4.2, 4.3 and 4.4 of
161 the Supplementary Material.

162 `McLust`

163 `McLust (52)` is a function from the R package `mcLust`. It estimates Gaussian mixture models for K clusters
164 based upon the maximum likelihood estimator of the parameters. It initialises upon a hierarchical clustering
165 of the data cut to K clusters. A range of choices of K and different covariance structures are compared and

Consensus clustering for Bayesian mixture models

166 the “best” selected using the Bayesian information criterion, (53) (details in section 4.2 of the Supplementary
167 Material).

168 **Bayesian inference**

169 To assess within-chain convergence of our Bayesian inference we use the Geweke Z -score statistic (54). Of
170 the chains that appear to behave properly we then assess across-chain convergence using \hat{R} (55) and the recent
171 extension provided by (56). If a chain has reached its stationary distribution the Geweke Z -score statistic is
172 expected to be normally distributed. Normality is tested for using a Shapiro-Wilks test (57). If a chain fails
173 this test (i.e., the associated p -value is less than 0.05), we assume that it has not achieved stationarity and it is
174 excluded from the remainder of the analysis. The samples from the remaining chains are then pooled and a
175 posterior similarity matrix (**PSM**) constructed. We use the `maxpear` function to infer a point clustering. For
176 more details see section 4.3 of the Supplementary Material.

177 **Analysis of the cell cycle in budding yeast**

178 **Datasets**

179 The cell cycle is crucial to biological growth, repair, reproduction, and development (58, 59, 60) and is highly
180 conserved among eukaryotes (60). . This means that understanding of the cell cycle of *S. cerevisiae* can
181 provide insight into a variety of cell cycle perturbations including those that occur in human cancer (61, 59)
182 and ageing (62). We aim to create clusters of genes that are co-expressed, have common regulatory proteins
183 and share a biological function. To achieve this, we use three datasets that were generated using different
184 omics technologies and target different aspects of the molecular biology underpinning the cell cycle process.

- 185 • Microarray profiles of RNA expression from (63), comprising measurements of cell-cycle-regulated
186 gene expression at 5-minute intervals for 200 minutes (up to three cell division cycles) and is referred
187 to as the **time course** dataset. The cells are synchronised at the START checkpoint in late G1-phase
188 using alpha factor arrest (63). We include only the genes identified by (63) as having periodic
189 expression profiles.
- 190 • Chromatin immunoprecipitation followed by microarray hybridization (**ChIP-chip**) data from (64).
191 This dataset discretizes p -values from tests of association between 117 DNA-binding transcriptional
192 regulators and a set of yeast genes. Based upon a significance threshold these p -values are represented
193 as either a 0 (no interaction) or a 1 (an interaction).

Consensus clustering for Bayesian mixture models

194 • Protein-protein interaction (**PPI**) data from BioGrid (65). This database consists of of physical
195 and genetic interactions between gene and gene products, with interactions either observed in high
196 throughput experiments or computationally inferred. The dataset we used contained 603 proteins
197 as columns. An entry of 1 in the $(i, j)^{th}$ cell indicates that the i^{th} gene has a protein product that is
198 believed to interact with the j^{th} protein.

199 The datasets were reduced to the 551 genes with no missing data in the PPI and ChIP-chip data, as in (35).

200 **Multiple dataset integration**

201 We applied consensus clustering to MDI for our integrative analysis. Details of MDI are in section 2.2 of the
202 Supplementary Material, but in short MDI jointly models the clustering in each dataset, inferring individual
203 clusterings for each dataset. These partitions are informed by similar structure in the other datasets, with MDI
204 learning this similarity as it models the partitions. The model does not assume global structure. This means
205 that the similarity between datasets is not strongly assumed in our model; individual clusters or genes that
206 align across datasets are based solely upon the evidence present in the data and not due to strong modelling
207 assumptions. Thus, datasets that share less common information can be included without fearing that this
208 will warp the final clusterings in some way.

209 The datasets were modelled using a mixture of Gaussian processes in the time course dataset and Multinomial
210 distributions in the ChIP-chip and PPI datasets.

211 **Results**

212 **Simulated data**

213 We use the ARI between the generating labels and the inferred clustering of each method to be our metric
214 of predictive performance. In figure 2, we see Mclust performs very well in the 2D and Small N , large P
215 scenarios, correctly identifying the true structure. However, the irrelevant features scenario sees a collapse in
216 performance, Mclust is blinded by the irrelevant features and identifies a clustering of $K = 1$.

217 The pooled samples from multiple long chains performs very well across all scenarios and appears to act as
218 an upper bound on the more practical implementations of consensus clustering.

219 Consensus clustering does uncover some of the generating structure in the data, even using a small number
220 of short chains. With sufficiently large ensembles and chain depth, consensus clustering is close to the
221 pooled Bayesian samples in predictive performance. It appears that for a constant chain depth increasing

Consensus clustering for Bayesian mixture models

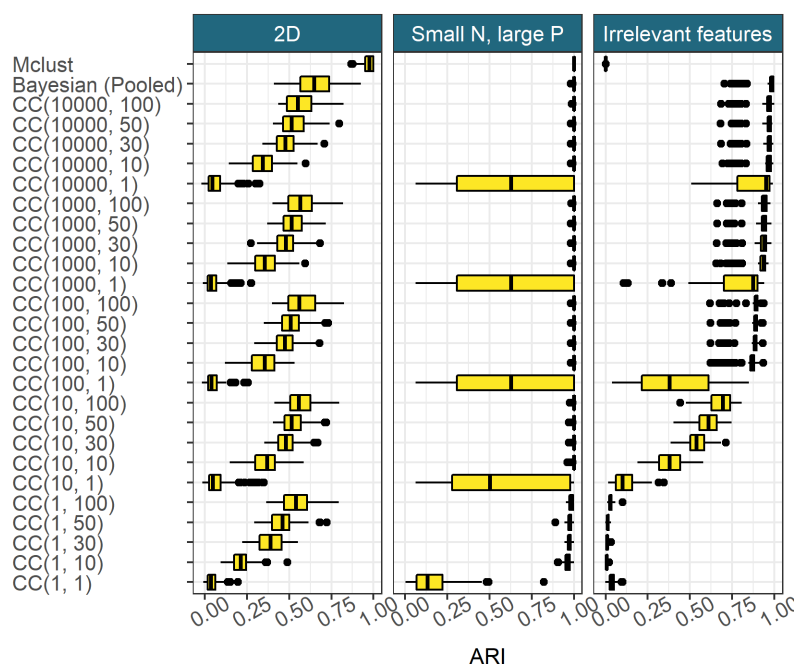


Figure 2: Model performance in the 100 simulated datasets for each scenario, defined as the ARI between the generating labels and the inferred clustering. $CC(d, w)$ denotes consensus clustering using the clustering from the d^{th} iteration from w different chains.

222 the ensemble width used follows a pattern of diminishing returns. There are strong initial gains for a greater
 223 ensemble width, but the improvement decreases for each successive chain. A similar pattern emerges in
 224 increasing chain length for a constant number of chains (figure 2).

225 For the PSMs from the individual chains, all entries are 0 or 1 (figure 3). This means only a single clustering
 226 is sampled within each chain, implying very little uncertainty in the partition. However, three different
 227 clustering solutions emerge across the chains, indicating that each individual chain is failing to explore the
 228 full support of the posterior distribution of the clustering. In general, while MCMC convergence theorems
 229 hold as the number of iterations tend to infinity, any finite chain might suffer in representing the full support
 230 of the posterior distribution, as we observe here. Moreover, the mixing of each chain can be poor as well (i.e.
 231 it may take a long time to reach the stationary distribution from an arbitrary initialisation). In our empirical
 232 study, we find that using many short runs provide similar point and interval estimates to running a small
 233 number of long chains (figure 3), while being computationally less expensive (figure 4), and hence more
 234 convenient for our applications.

235 Figure 4 shows that chain length is directly proportional to the time taken for the chain to run. This means
 236 that using an ensemble of shorter chains, as in consensus clustering, can offer large reductions in the time
 237 cost of analysis when a parallel environment is available compared to standard Bayesian inference. Even on a

Consensus clustering for Bayesian mixture models

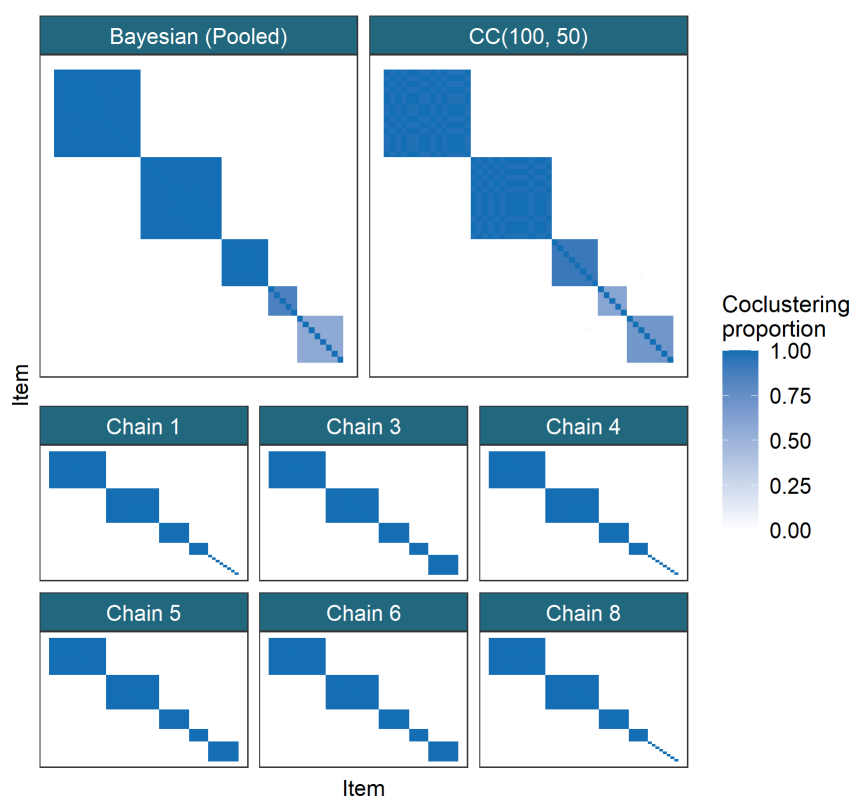


Figure 3: Comparison of similarity matrices from a dataset for the Small N , large P scenario. In each matrix, the $(i, j)^{th}$ entry is the proportion of clusterings for which the i^{th} and j^{th} items co-clustered for the method in question. In the first row the PSM of the pooled Bayesian samples is compared to the CM for CC(100, 50), with a common ordering of rows and columns in both heatmaps. In the following rows, 6 of the long chains that passed the tests of convergence are shown.

238 laptop of 8 cores running an ensemble of 1,000 chains of length 1,000 will require approximately half as
239 much time as running 10 chains of length 100,000 due to parallelisation, and the potential benefits are far
240 greater when using a large computing cluster.

241 Additional results for these and other simulations are in section 4.4 of the Supplementary Material.

242 **Multi-omics analysis of the cell cycle in budding yeast**

243 We use the stopping rule proposed in to determine our ensemble depth and width. In figure 5, we see that the
244 change in the consensus matrices from increasing the ensemble depth and width is diminishing in keeping
245 with results in the simulations. We see no strong improvement after $D = 6,000$ and increasing the number
246 of learners from 500 to 1,000 has small effect. We therefore use the largest ensemble available, a depth

Consensus clustering for Bayesian mixture models

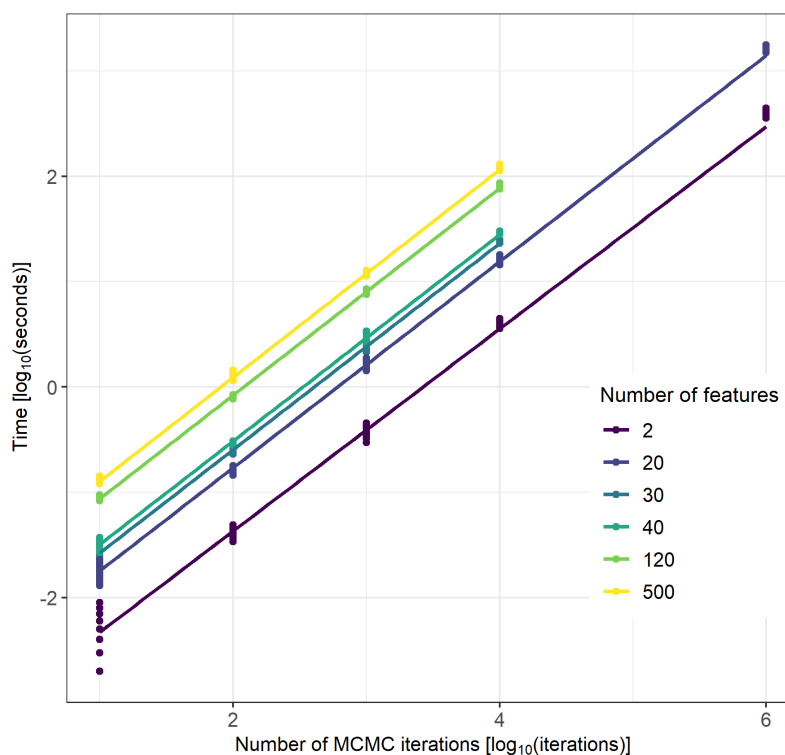


Figure 4: The time taken for different numbers of iterations of MCMC moves in $\log_{10}(\text{seconds})$. The relationship between chain length, D , and the time taken is linear (the slope is approximately 1 on the \log_{10} scale), with a change of intercept, for different dimensions. The runtime of each Markov chain was recorded using the terminal command `time`, measured in milliseconds.

247 $D = 10001$ and width $W = 1000$, believing this ensemble is stable (additional evidence in section 5.1 of the
248 Supplementary Material).

249 We focus upon the genes that tend to have the same cluster label across multiple datasets. More formally, we
250 analyse the clustering structure among genes for which $\hat{P}(c_{nl} = c_{nm}) > 0.5$, where c_{nl} denotes the cluster
251 label of gene n in dataset l . In our analysis it is the signal shared across the time course and ChIP-chip
252 datasets that is strongest, with 261 genes (nearly half of the genes present) in this pairing tending to have a
253 common label, whereas only 56 genes have a common label across all three datasets. Thus, we focus upon
254 this pairing of datasets in the results of the analysis performed using all three datasets. We show the gene
255 expression and regulatory proteins of these genes separated by their cluster in figure 6. In figure 6, the clusters
256 in the time series data have tight, unique signatures (having different periods, amplitudes, or both) and in the
257 ChIP-chip data clusters are defined by a small number of well-studied transcription factors (TFs) (see table 2
258 of the Supplementary Material for details of these TFs, many of which are well known to regulate cell cycle
259 expression, 66).

Consensus clustering for Bayesian mixture models

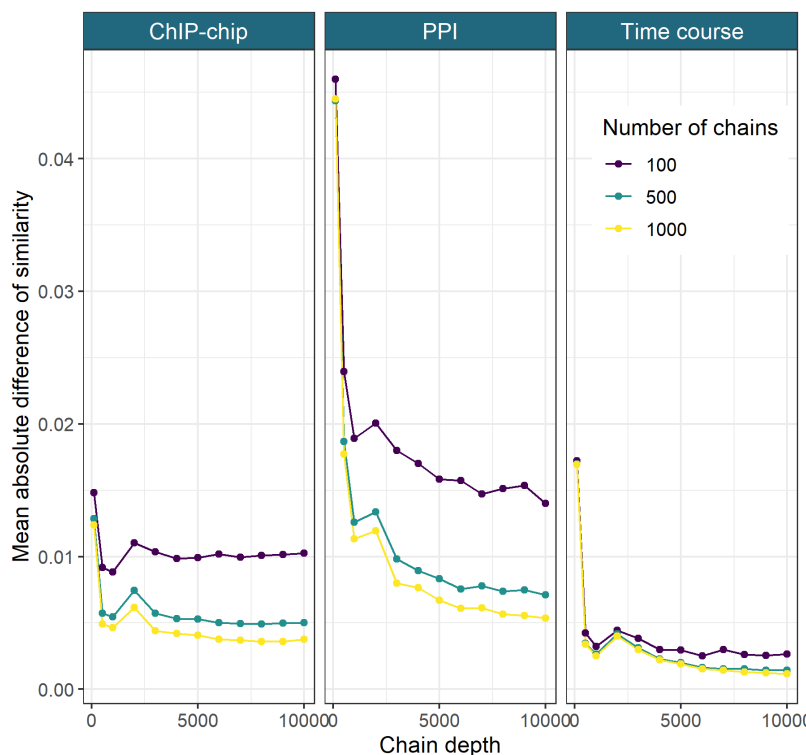


Figure 5: The mean absolute difference between the sequential Consensus matrices. For a set of chain lengths, $D' = \{d_1, \dots, d_I\}$ and number of chains, $W' = \{w_1, \dots, w_J\}$, we take the mean of the absolute difference between the consensus matrix for (d_i, w_j) and (d_{i-1}, w_j) (here $D' = \{101, 501, 1001, 2001, \dots, 10001\}$ and $W' = \{100, 500, 1000\}$).

260 As an example, we briefly analyse clusters 9 and 16 in greater depth. Cluster 9 has strong association with
261 MBP1 and some interactions with SWI6, as can be seen in figure 6. The Mbp1-Swi6p complex, MBF, is
262 associated with DNA replication (67). The first time point, 0 minutes, in the time course data is at the START
263 checkpoint, or the G1/S transition. The members of cluster 9 begin highly expressed at this point before
264 quickly dropping in expression (in the first of the 3 cell cycles). This suggests that many transcripts are
265 produced immediately in advance of S-phase, and thus are required for the first stages of DNA synthesis.
266 These genes' descriptions (found using `org.Sc.sgd.db`, 68, and shown in table 3 of the Supplementary
267 Material) support this hypothesis, as many of the members are associated with DNA replication, repair and/or
268 recombination. Additionally, *TOF1*, *MRC1* and *RAD53*, members of the replication checkpoint (69, 70)
269 emerge in the cluster as do members of the cohesin complex. Cohesin is associated with sister chromatid
270 cohesion which is established during the S-phase of the cell cycle (71) and also contributes to transcription
271 regulation, DNA repair, chromosome condensation, homolog pairing (72), fitting the theme of cluster 9.
272 Cluster 16 appears to be a cluster of S-phase genes, consisting of *GAS3*, *NRM1* and *PDS1* and the genes
273 encoding the histones H1, H2A, H2B, H3 and H4. Histones are the chief protein components of chromatin

Consensus clustering for Bayesian mixture models

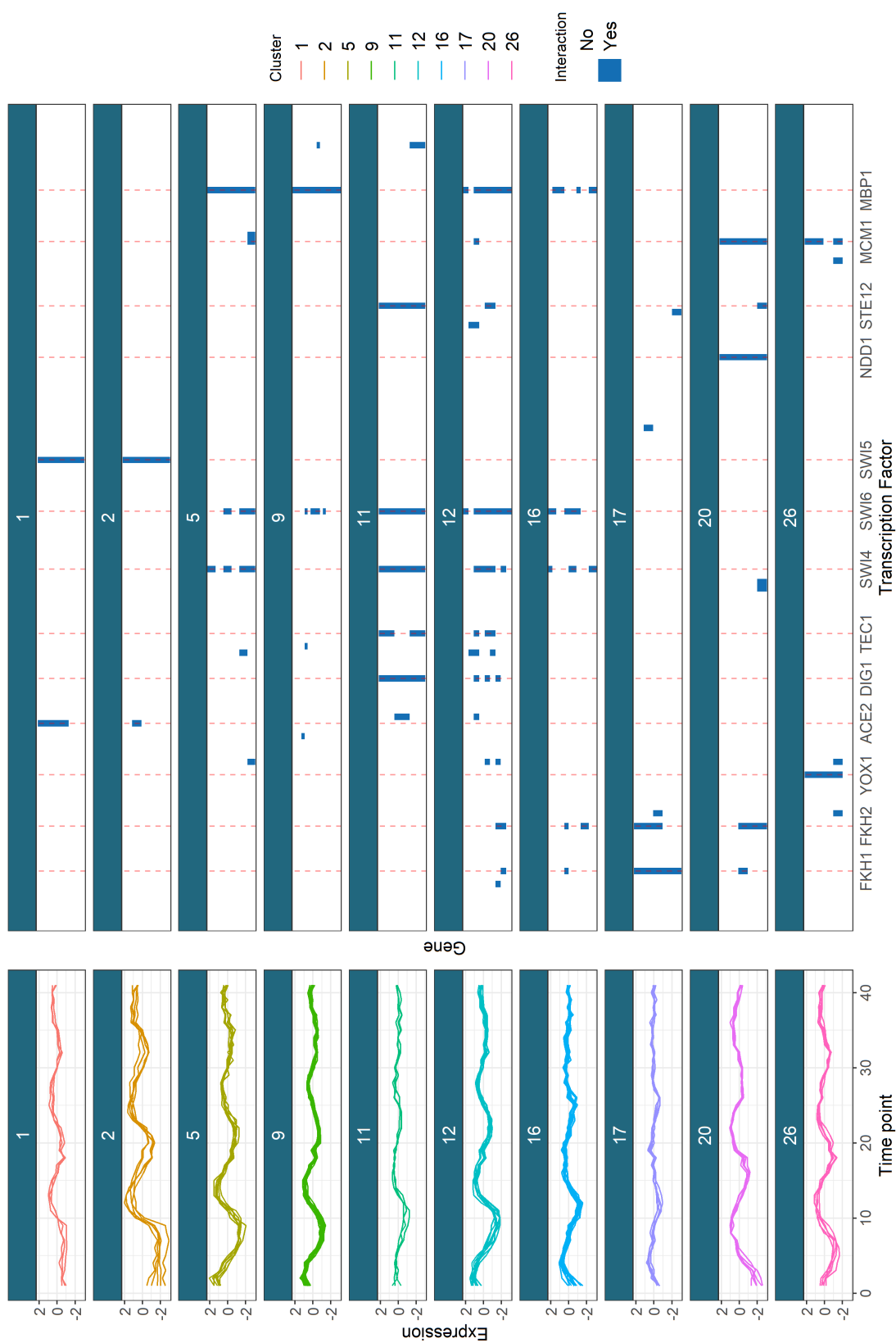


Figure 6: The gene clusters which tend to have a common label across the time course and ChIP-chip datasets, shown in these datasets. We include only the clusters with more than one member and more than half the members having some interactions in the ChIP-chip data. Red lines for the most common transcription factors are included.

Consensus clustering for Bayesian mixture models

274 (73) and are important contributors to gene regulation (74). They are known to peak in expression in
275 S-phase (63), which matches the first peak of this cluster early in the time series. Of the other members,
276 *NRMI* is a transcriptional co-repressor of MBF-regulated gene expression acting at the transition from G1
277 to S-phase (75, 76). Pds1p binds to and inhibits the Esp1 class of sister separating proteins, preventing
278 sister chromatids separation before M-phase (77, 71). *GAS3*, is not well studied. It interacts with *SMT3*
279 which regulates chromatid cohesion, chromosome segregation and DNA replication (among other things).
280 Chromatid cohesion ensures the faithful segregation of chromosomes in mitosis and in both meiotic divisions
281 (78) and is instantiated in S-phase (71). These results, along with the very similar expression profile to the
282 histone genes in the time course data, suggest that *GAS3* may be more directly involved in DNA replication
283 or chromatid cohesion than is currently believed.

284 We attempt to perform a similar analysis using traditional Bayesian inference of MDI, but after 36 hours of
285 runtime there is no consistency or convergence across chains. We use the Geweke statistic and \hat{R} to reduce to
286 the five best behaved chains (none of which appear to be converged, see section 5.2 of the Supplementary
287 Material for details). If we then compare the distribution of sampled values for the ϕ parameters for these
288 long chains, the final ensemble used ($D = 10001$, $W = 1000$) and the pooled samples from the 5 long chains,
289 then we see that the distribution of the pooled samples from the long chains (which might be believed to
290 sampling different parts of the posterior distribution) is closer in appearance to the distributions sampled by
291 the consensus clustering than to any single chain (figure 7). Further disagreement between chains is shown in
292 the Gene Ontology term over-representation analysis in section 5.3 of the Supplementary Material.

293 Discussion

294 Our proposed method has demonstrated good performance on simulation studies, uncovering the generating
295 structure in many cases and performing comparably to Mclust and long chains in many scenarios. We saw
296 that when the chains are sufficiently deep that the ensemble approximates Bayesian inference, as shown by
297 the similarity between the PSMs and the CM in the 2D scenario where the individual chains do not become
298 trapped in a single mode. We have shown cases where many short runs are computationally less expensive
299 than one long chain and give meaningful point and interval estimates; estimates that are very similar to those
300 from the limiting case of a Markov chain. Thus if individual chains are suffering from mixing problems or
301 are too computationally expensive to run, consensus clustering may provide a viable option. We also showed
302 that the ensemble of short chains is more robust to irrelevant features than Mclust.

Consensus clustering for Bayesian mixture models

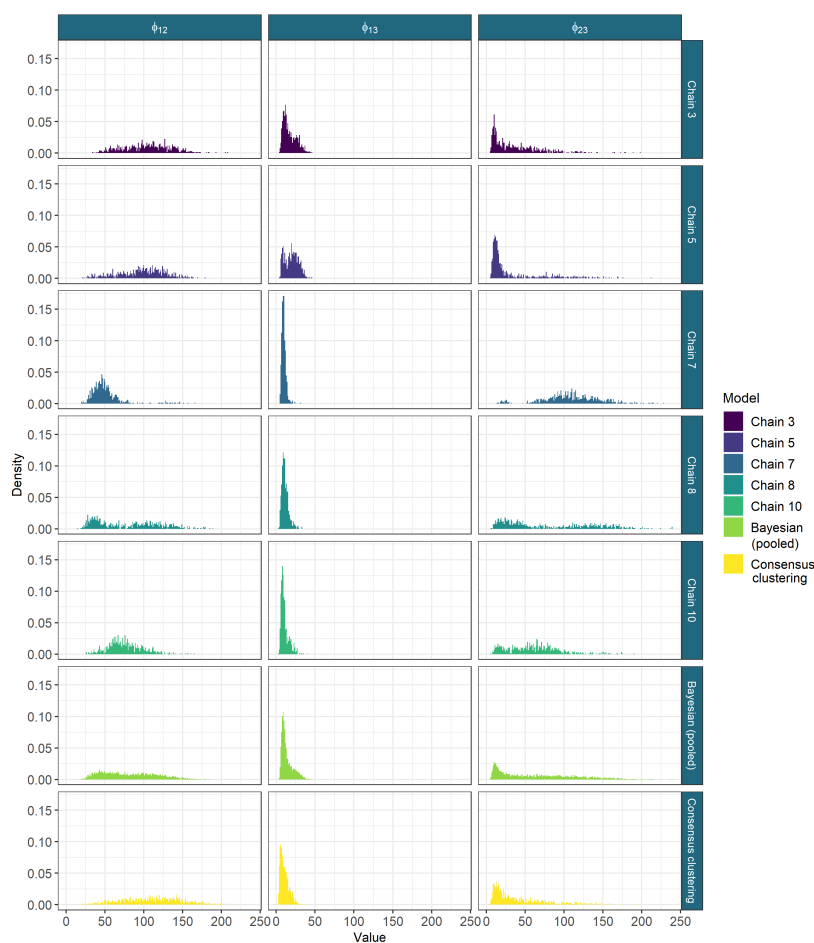


Figure 7: The sampled values for the ϕ parameters from the long chains, their pooled samples and the consensus using 1000 chains of depth 10,001. The long chains display a variety of behaviours. Across chains there is no clear consensus on the nature of the posterior distribution. The samples from any single chain are not particularly close to the behaviour of the pooled samples across all three parameters. It is the consensus clustering that most approaches this pooled behaviour.

303 We proposed a method of assessing ensemble stability and deciding upon ensemble size which we used when
304 performing an integrative analysis of yeast cell cycle data using MDI, an extension of Bayesian mixture
305 models that jointly models multiple datasets. We uncovered many genes with shared signal across several
306 datasets and explored the meaning of some of the inferred clusters using data external to the analysis. We
307 found biologically meaningful results as well as signal for possibly novel biology. We also showed that
308 individual chains for the existing implementation of MDI do not converge in a practical length of time, having
309 run 10 chains for 36 hours with no consistent behaviour across chains. This means that Bayesian inference of
310 the MDI model is not practical on this dataset with the software currently available.

311 However, consensus clustering does lose the theoretical framework of true Bayesian inference. We attempt to
312 mitigate this with our assessment of stability in the ensemble, but this diagnosis is heuristic and subjective, and

Consensus clustering for Bayesian mixture models

313 while there is empirical evidence for its success, it lacks the formal results for the tests of model convergence
314 for Bayesian inference.

315 More generally, we have benchmarked the use of an ensemble of Bayesian mixture models, showing that this
316 approach can infer meaningful clusterings and overcomes the problem of multi-modality in the likelihood
317 surface even in high dimensions, thereby providing more stable clusterings than individual long chains that
318 are prone to becoming trapped in individual modes. We also show that the ensemble can be significantly
319 quicker to run. In our multi-omics study we have demonstrated that the method can be applied as a wrapper to
320 more complex Bayesian clustering methods using existing implementations and that this provides meaningful
321 results even when individual chains fail to converge. This enables greater application of complex Bayesian
322 clustering methods without requiring re-implementation using more clever MCMC methods, a process that
323 would involve a significant investment of human time.

324 We expect that researchers interested in applying some of the Bayesian integrative clustering models such
325 as MDI and Clusternomics (36) will be enabled to do so, as consensus clustering overcomes some of the
326 unwieldiness of existing implementations of these complex models. More generally, we expect that our
327 method will be useful to researchers performing cluster analysis of high-dimensional data where the runtime
328 of MCMC methods becomes too onerous and multi-modality is more likely to be present.

329 **Abbreviations**

330 ARI: Adjusted Rand Index

331 ChIP-chip: Chromatin immunoprecipitation followed by microarray hybridization

332 CM: Consensus Matrix

333 MCMC: Markov chain Monte Carlo

334 MDI: Multiple Dataset Integration

335 PCA: Principal Component Analysis

336 PPI: Protein-Protein Interaction

337 PSM: Posterior Similarity Matrix

338 SSE: Sum of Squared Errors

339 TF: Transcription Factor

340 **Availability of data and materials**

341 The code and datasets supporting the conclusions of this article are available in the github repository,
342 <https://github.com/stcolema/ConsensusClusteringForBayesianMixtureModels>.

343 **Competing interests**

344 The authors declare that they have no competing interests.

345 **Funding**

346 This work was funded by the MRC (MC UU 00002/4, MC UU 00002/13) and supported by the NIHR
347 Cambridge Biomedical Research Centre (BRC-1215-20014). The views expressed are those of the author(s)
348 and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care. This research
349 was funded in whole, or in part, by the Wellcome Trust [WT107881]. For the purpose of Open Access, the
350 author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising
351 from this submission.

352 **Authors' contributions**

353 SC designed the simulation study with contributions from PK and CW, performed the analyses and wrote
354 the manuscript. PK and CW provided an equal contribution of joint supervision, directing the research and
355 provided suggestions such as the stopping rule. All contributed to interpreting the results of the analyses. All
356 authors revised and approved the final manuscript.

357 **References**

- 358 [1] Hejblum BP, Skinner J, Thiébaud R. Time-course gene set analysis for longitudinal gene expression
359 data. *PLoS computational biology*. 2015;11(6):e1004310.
- 360 [2] Bai JP, Alekseyenko AV, Statnikov A, Wang IM, Wong PH. Strategic applications of gene expression:
361 from drug discovery/development to bedside. *The AAPS journal*. 2013;15(2):427–437.
- 362 [3] Emmert-Streib F, Dehmer M, Haibe-Kains B. Gene regulatory networks and their applications: under-
363 standing biological and medical problems in terms of networks. *Frontiers in cell and developmental*
364 *biology*. 2014;2:38.

Consensus clustering for Bayesian mixture models

- 365 [4] Lloyd S. Least squares quantization in PCM. *IEEE transactions on information theory*. 1982;28(2):129–
366 137.
- 367 [5] Forgy EW. Cluster analysis of multivariate data: efficiency versus interpretability of classifications.
368 *biometrics*. 1965;21:768–769.
- 369 [6] Rousseeuw PJ. Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis.
370 *Journal of Computational and Applied Mathematics*. 1987 Nov;20:53–65.
- 371 [7] Arthur D, Vassilvitskii S. *k-means++: The advantages of careful seeding*. Stanford; 2006.
- 372 [8] Breiman L. Random forests. *Machine learning*. 2001;45(1):5–32.
- 373 [9] Friedman JH. Stochastic gradient boosting. *Computational statistics & data analysis*. 2002;38(4):367–
374 378.
- 375 [10] Monti S, Tamayo P, Mesirov J, Golub T. Consensus clustering: a resampling-based method for class
376 discovery and visualization of gene expression microarray data. *Machine learning*. 2003;52(1-2):91–
377 118.
- 378 [11] Wilkerson, D M, Hayes, Neil D. ConsensusClusterPlus: a class discovery tool with confidence
379 assessments and item tracking. *Bioinformatics*. 2010;26(12):1572–1573.
- 380 [12] John CR, Watson D, Russ D, Goldmann K, Ehrenstein M, Pitzalis C, et al. M3C: Monte Carlo
381 reference-based consensus clustering. *Scientific reports*. 2020;10(1):1–14.
- 382 [13] Gu Z, Schlesner M, Hübschmann D. cola: an R/Bioconductor package for consensus partitioning
383 through a general framework. *Nucleic Acids Research*. 2020 12;Gkaa1146. Available from: <https://doi.org/10.1093/nar/gkaa1146>.
384
- 385 [14] Lehmann BD, Bauer JA, Chen X, Sanders ME, Chakravarthy AB, Shyr Y, et al. Identification of human
386 triple-negative breast cancer subtypes and preclinical models for selection of targeted therapies. *The*
387 *Journal of clinical investigation*. 2011;121(7):2750–2767.
- 388 [15] Verhaak RG, Hoadley KA, Purdom E, Wang V, Qi Y, Wilkerson MD, et al. Integrated genomic analysis
389 identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA,
390 IDH1, EGFR, and NF1. *Cancer cell*. 2010;17(1):98–110.
- 391 [16] Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3: consensus clustering
392 of single-cell RNA-seq data. *Nature methods*. 2017;14(5):483–486.
- 393 [17] Li T, Ding C. Weighted Consensus Clustering. In: *Proceedings of the 2008 SIAM International*
394 *Conference on Data Mining*. Society for Industrial and Applied Mathematics; 2008. p. 798–809.

Consensus clustering for Bayesian mixture models

- 395 [18] Carpineto C, Romano G. Consensus Clustering Based on a New Probabilistic Rand Index with
396 Application to Subtopic Retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
397 2012 Dec;34(12):2315–2326.
- 398 [19] Strehl A, Ghosh J. Cluster Ensembles – A Knowledge Reuse Framework for Combining Multiple
399 Partitions. *Journal of Machine Learning Research*. 2002;3:583–617.
- 400 [20] Ghaemi R, Sulaiman MN, Ibrahim H, Mustapha N, et al. A survey: clustering ensembles techniques.
401 *World Academy of Science, Engineering and Technology*. 2009;50:636–645.
- 402 [21] Ünlü R, Xanthopoulos P. Estimating the Number of Clusters in a Dataset via Consensus Clustering.
403 *Expert Systems with Applications*. 2019 Jul;125:33–39.
- 404 [22] Fraley C, Raftery AE. Model-based clustering, discriminant analysis, and density estimation. *Journal*
405 *of the American statistical Association*. 2002;97(458):611–631.
- 406 [23] Fraley C. How Many Clusters? Which Clustering Method? Answers Via Model-Based Cluster Analysis.
407 *The Computer Journal*. 1998 Aug;41(8):578–588.
- 408 [24] Antoniak CE. Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems.
409 *The Annals of Statistics*. 1974 Nov;2(6):1152–1174.
- 410 [25] Ferguson TS. Bayesian Density Estimation by Mixtures of Normal Distributions. In: Rizvi MH, Rustagi
411 JS, Siegmund D, editors. *Recent Advances in Statistics*. Academic Press; 1983. p. 287–302.
- 412 [26] Lo AY. On a Class of Bayesian Nonparametric Estimates: I. Density Estimates. *The Annals of Statistics*.
413 1984;12(1):351–357.
- 414 [27] Richardson S, Green PJ. On Bayesian analysis of mixtures with an unknown number of components.
415 *Journal of the Royal Statistical Society: series B*. 1997;59(4):731–792.
- 416 [28] Miller JW, Harrison MT. Mixture models with a prior on the number of components. *Journal of the*
417 *American Statistical Association*. 2018;113(521):340–356.
- 418 [29] Rousseau J, Mengersen K. Asymptotic behaviour of the posterior distribution in overfitted mixture
419 models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2011;73(5):689–
420 710.
- 421 [30] Medvedovic M, Sivaganesan S. Bayesian infinite mixture model based clustering of gene expression
422 profiles. *Bioinformatics*. 2002;18(9):1194–1206.
- 423 [31] Chan C, Feng F, Ottinger J, Foster D, West M, Kepler TB. Statistical mixture modeling for cell
424 subtype identification in flow cytometry. *Cytometry Part A: The Journal of the International Society for*
425 *Analytical Cytology*. 2008;73(8):693–701.

Consensus clustering for Bayesian mixture models

- 426 [32] Hejblum BP, Alkassim C, Gottardo R, Caron F, Thiébaud R, et al. Sequential Dirichlet process mixtures
427 of multivariate skew t -distributions for model-based clustering of flow cytometry data. *The Annals of*
428 *Applied Statistics*. 2019;13(1):638–660.
- 429 [33] Prabhakaran S, Azizi E, Carr A, Pe'er D. Dirichlet process mixture model for correcting technical
430 variation in single-cell gene expression data. In: *International Conference on Machine Learning*; 2016.
431 p. 1070–1079.
- 432 [34] Crook OM, Mulvey CM, Kirk PD, Lilley KS, Gatto L. A Bayesian mixture modelling approach for
433 spatial proteomics. *PLoS computational biology*. 2018;14(11):e1006516.
- 434 [35] Kirk P, Griffin JE, Savage RS, Ghahramani Z, Wild DL. Bayesian correlated clustering to integrate
435 multiple datasets. *Bioinformatics*. 2012;28(24):3290–3297.
- 436 [36] Gabasova E, Reid J, Wernisch L. Clusternomics: Integrative context-dependent clustering for heteroge-
437 neous datasets. *PLoS computational biology*. 2017;13(10):e1005781.
- 438 [37] Strauss ME, Kirk PD, Reid JE, Wernisch L. GPseudoClust: deconvolution of shared pseudo-profiles at
439 single-cell resolution. *Bioinformatics*. 2020;36(5):1484–1491.
- 440 [38] Robert CP, Elvira V, Tawn N, Wu C. Accelerating MCMC algorithms. *Wiley Interdisciplinary Reviews:*
441 *Computational Statistics*. 2018;10(5):e1435.
- 442 [39] Neiswanger W, Wang C, Xing EP. Asymptotically Exact, Embarrassingly Parallel MCMC. In:
443 *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence. UAI'14. Arlington,*
444 *Virginia, USA: AUAI Press; 2014. p. 623–632.*
- 445 [40] Murray L. Distributed Markov Chain Monte Carlo. In: *Proceedings of Neural Information Processing*
446 *Systems workshop on learning on cores, clusters and clouds. vol. 11; 2010. .*
- 447 [41] Jacob PE, O'Leary J, Atchadé YF. Unbiased Markov Chain Monte Carlo Methods with Couplings.
448 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2020;82(3):543–600.
- 449 [42] Wold S, Esbensen K, Geladi P. Principal component analysis. *Chemometrics and intelligent laboratory*
450 *systems*. 1987;2(1-3):37–52.
- 451 [43] Fritsch A, Ickstadt K, et al. Improved criteria for clustering based on the posterior similarity matrix.
452 *Bayesian analysis*. 2009;4(2):367–391.
- 453 [44] Fritsch A. mcclust: process an MCMC sample of clusterings; 2012. R package version 1.0. Available
454 from: <https://CRAN.R-project.org/package=mcclust>.
- 455 [45] Wade S, Ghahramani Z. Bayesian Cluster Analysis: Point Estimation and Credible Balls (with
456 Discussion). *Bayesian Analysis*. 2018 Jun;13(2):559–626.

Consensus clustering for Bayesian mixture models

- 457 [46] Lourenço A, Rota Bulò S, Rebagliati N, Fred ALN, Figueiredo MAT, Pelillo M. Probabilistic Consensus
458 Clustering Using Evidence Accumulation. *Machine Learning*. 2015 Jan;98(1):331–357.
- 459 [47] Dahl DB, Johnson DJ, Mueller P. Search Algorithms and Loss Functions for Bayesian Clustering.
460 arXiv:210504451 [stat]. 2021 May;.
- 461 [48] Von Luxburg U, Ben-David S. Towards a statistical theory of clustering. In: *Pascal workshop on
462 statistics and optimization of clustering*. Citeseer; 2005. p. 20–26.
- 463 [49] Meinshausen N, Bühlmann P. Stability selection. *Journal of the Royal Statistical Society: Series B
464 (Statistical Methodology)*. 2010;72(4):417–473.
- 465 [50] Law MH, Jain AK, Figueiredo M. Feature selection in mixture-based clustering. In: *Advances in neural
466 information processing systems*; 2003. p. 641–648.
- 467 [51] Hubert L, Arabie P. Comparing partitions. *Journal of classification*. 1985;2(1):193–218.
- 468 [52] Scrucca L, Fop M, Murphy BT, Raftery AE. mclust 5: clustering, classification and density estimation
469 using Gaussian finite mixture models. *The R Journal*. 2016;8(1):289–317. Available from: <https://doi.org/10.32614/RJ-2016-021>.
- 470
- 471 [53] Schwarz G, et al. Estimating the dimension of a model. *The annals of statistics*. 1978;6(2):461–464.
- 472 [54] Geweke J, et al. Evaluating the accuracy of sampling-based approaches to the calculation of posterior
473 moments. vol. 196. Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN;
474 1991.
- 475 [55] Gelman A, Rubin DB, et al. Inference from iterative simulation using multiple sequences. *Statistical
476 science*. 1992;7(4):457–472.
- 477 [56] Vats D, Knudson C. Revisiting the Gelman-Rubin diagnostic. arXiv preprint arXiv:181209384. 2018;.
- 478 [57] Shapiro SS, Wilk MB. An analysis of variance test for normality (complete samples). *Biometrika*.
479 1965;52(3/4):591–611.
- 480 [58] Tyson JJ, Chen KC, Novák B. Cell Cycle, Budding Yeast. In: Dubitzky W, Wolkenhauer O, Cho KH,
481 Yokota H, editors. *Encyclopedia of Systems Biology*. New York, NY: Springer New York; 2013. p.
482 337–341.
- 483 [59] Chen KC, Calzone L, Csikasz-Nagy A, Cross FR, Novak B, Tyson JJ. Integrative analysis of cell cycle
484 control in budding yeast. *Molecular biology of the cell*. 2004;15(8):3841–3862.
- 485 [60] Alberts B, Johnson A, Lewis J, Raff M, Roberts K, Walter P. The cell cycle and programmed cell death.
486 *Molecular biology of the cell*. 2002;4:983–1027.

Consensus clustering for Bayesian mixture models

- 487 [61] Ingalls B, Duncker B, Kim D, McConkey B. Systems level modeling of the cell cycle using budding
488 yeast. *Cancer informatics*. 2007;3:117693510700300020.
- 489 [62] Jiménez J, Bru S, Ribeiro M, Clotet J. Live fast, die soon: cell cycle progression and lifespan in yeast
490 cells. *Microbial Cell*. 2015;2(3):62.
- 491 [63] Granovskaia MV, Jensen LJ, Ritchie ME, Toedling J, Ning Y, Bork P, et al. High-resolution transcription
492 atlas of the mitotic cell cycle in budding yeast. *Genome biology*. 2010;11(3):1–11.
- 493 [64] Harbison CT, Gordon DB, Lee TI, Rinaldi NJ, Macisaac KD, Danford TW, et al. Transcriptional
494 regulatory code of a eukaryotic genome. *Nature*. 2004;431(7004):99–104.
- 495 [65] Stark C, Breitkreutz BJ, Reguly T, Boucher L, Breitkreutz A, Tyers M. BioGRID: a general repository
496 for interaction datasets. *Nucleic acids research*. 2006;34(suppl_1):D535–D539.
- 497 [66] Simon I, Barnett J, Hannett N, Harbison CT, Rinaldi NJ, Volkert TL, et al. Serial regulation of
498 transcriptional regulators in the yeast cell cycle. *Cell*. 2001;106(6):697–708.
- 499 [67] Iyer VR, Horak CE, Scafe CS, Botstein D, Snyder M, Brown PO. Genomic binding sites of the yeast
500 cell-cycle transcription factors SBF and MBF. *Nature*. 2001;409(6819):533–538.
- 501 [68] Carlson M, Falcon S, Pages H, Li N. Org. sc. sgd. db: Genome wide annotation for yeast. R package
502 version. 2014;2(1).
- 503 [69] Bando M, Katou Y, Komata M, Tanaka H, Itoh T, Sutani T, et al. Csm3, Tof1, and Mrc1 form a
504 heterotrimeric mediator complex that associates with DNA replication forks. *Journal of Biological
505 Chemistry*. 2009;284(49):34355–34365.
- 506 [70] Lao JP, Ulrich KM, Johnson JR, Newton BW, Vashisht AA, Wohlschlegel JA, et al. The yeast DNA
507 damage checkpoint kinase Rad53 targets the exoribonuclease, Xrn1. *G3: Genes, Genomes, Genetics*.
508 2018;8(12):3931–3944.
- 509 [71] Tóth A, Ciosk R, Uhlmann F, Galova M, Schleiffer A, Nasmyth K. Yeast cohesin complex requires
510 a conserved protein, Eco1p (Ctf7), to establish cohesion between sister chromatids during DNA
511 replication. *Genes & development*. 1999;13(3):320–333.
- 512 [72] Mehta GD, Kumar R, Srivastava S, Ghosh SK. Cohesin: functions beyond sister chromatid cohesion.
513 *FEBS letters*. 2013;587(15):2299–2312.
- 514 [73] Fischle W, Wang Y, Allis CD. Histone and chromatin cross-talk. *Current opinion in cell biology*.
515 2003;15(2):172–183.
- 516 [74] Bannister AJ, Kouzarides T. Regulation of chromatin by histone modifications. *Cell research*.
517 2011;21(3):381–395.

Consensus clustering for Bayesian mixture models

- 518 [75] de Bruin RA, Kalashnikova TI, Chahwan C, McDonald WH, Wohlschlegel J, Yates III J, et al. Con-
519 straining G1-specific transcription to late G1 phase: the MBF-associated corepressor Nrm1 acts via
520 negative feedback. *Molecular cell*. 2006;23(4):483–496.
- 521 [76] Aligianni S, Lackner DH, Klier S, Rustici G, Wilhelm BT, Marguerat S, et al. The fission yeast
522 homeodomain protein Yox1p binds to MBF and confines MBF-dependent cell-cycle transcription to
523 G1-S via negative feedback. *PLoS Genet*. 2009;5(8):e1000626.
- 524 [77] Ciosk R, Zachariae W, Michaelis C, Shevchenko A, Mann M, Nasmyth K. An ESP1/PDS1 complex
525 regulates loss of sister chromatid cohesion at the metaphase to anaphase transition in yeast. *Cell*.
526 1998;93(6):1067–1076.
- 527 [78] Cooper KF, Mallory MJ, Guacci V, Lowe K, Strich R. Pds1p is required for meiotic recombination and
528 prophase I progression in *Saccharomyces cerevisiae*. *Genetics*. 2009;181(1):65–79.

529 **Additional Files**

530 **Additional file 1 — Supplementary materials**

531 Additional relevant theory, background and results. This includes some more formal definitions, details of
532 Bayesian mixture models and MDI, the general consensus clustering algorithm, additional simulations and
533 the generating algorithm used, steps in assessing Bayesian model convergence in both the simulated datasets
534 and yeast analysis, a table of the transcription factors that define the clustering in the ChIP-chip dataset, a
535 table of the gene descriptions for some of the clusters that emerge across the time course and ChIP-chip
536 datasets and Gene Ontology term over-representation analysis of the clusterings from the yeast datasets.
537 (PDF, 10MB)

Consensus clustering for Bayesian mixture models: Supplementary materials

Stephen Coleman, Paul DW Kirk and Chris Wallace

1 Definitions

Definition 1 (Coclustering matrix) *The coclustering matrix describes a clustering or partition in a binary matrix with the $(i, j)^{\text{th}}$ entry indicating if items i and j are allocated to the same cluster.*

Definition 2 (Consensus matrix) *Given W clusterings for a dataset of N items, $c_s = (c_{s1}, \dots, c_{sN})$, the consensus matrix is a $N \times N$ matrix where the $(i, j)^{\text{th}}$ entry records the proportions of clusterings for which items i and j are allocated the same label. More formally, it is the matrix \mathbb{C} such that*

$$\mathbb{C}(i, j) = \frac{1}{W} \sum_{s=1}^W \mathbf{I}(c_{si} = c_{sj}) \quad (1)$$

where $\mathbb{I}(\cdot)$ is the indicator function taking a value of 1 if the argument is true and 0 otherwise.

Definition 3 (Posterior similarity matrix) *A consensus matrix for which all the clusterings are generated from a converged Markov chain for some Bayesian clustering model. Sometimes abbreviated to PSM.*

Definition 4 (Partition or Clustering) *For a dataset of items $X = (x_1, \dots, x_N)$, a partition or clustering is a set of disjoint sets covering X , normally indicated by a N -vector of integers indicating which set each item is associated with. Note that these labels only have meaning relative to each other, they are symbolic. Each set within the clustering is referred to as a cluster.*

2 The models

2.1 Individual dataset

In the simulations (see section 4) where individual datasets are modelled a *Bayesian mixture model* is used. We write the basic mixture model for inde-

pendent items $X = (x_1, \dots, x_N)$ as

$$x_n \sim \sum_{k=1}^K \pi_k f(x_n | \theta_k) \quad \text{independently for } n = 1, \dots, N \quad (2)$$

where $f(\cdot | \theta)$ is some family of densities parametrised by θ . A common choice is the Gaussian density function, with $\theta = (\mu, \sigma^2)$ (as in our simulation study). K , the number of subgroups in the population, $\{\theta_k\}_{k=1}^K$, the component parameters, and $\pi = (\pi_1, \dots, \pi_K)$, the component weights are the objects to be inferred. In the context of *clustering*, such a model arises due to the belief that the population from which the random sample under analysis has been drawn consists of K unknown groups proportional to π . In this setting it is natural to include a latent *allocation variable*, $c = (c_1, \dots, c_N)$, to indicate which group each item is drawn from, with each non-empty component of the mixture corresponds to a cluster. The model is

$$\begin{aligned} p(c_n = k) &= \pi_k \quad \text{for } k = 1, \dots, K, \\ x_n | c_n &\sim f(x_n | \theta_k) \quad \text{independently for } n = 1, \dots, N. \end{aligned} \quad (3)$$

The joint model can then be written

$$p(X, c, K, \pi, \theta) = p(X | c, \pi, K, \theta) p(\theta | c, \pi, K) p(c | \pi, K) p(\pi | K) p(K)$$

We assume conditional independence between certain parameters such that the model reduces to

$$p(X, c, \theta, \pi, K) = p(\pi | K) p(\theta | K) p(K) \prod_{n=1}^N p(x_n | c_n, \theta_{c_n}) p(c_n | \pi, K). \quad (4)$$

Additional flexibility is provided by the inclusion of hyperparameters on the priors for π and θ , denoted α and η respectively. In our context where $\theta = (\mu, \sigma^2)$, we use

$$\sigma^2 \sim \Gamma^{-1}(a, b), \quad (5)$$

$$\mu \sim \mathcal{N}\left(\xi, \frac{1}{\lambda} \sigma^2\right), \quad (6)$$

$$\pi \sim \text{Dirichlet}(\alpha). \quad (7)$$

The directed acyclic graph (**DAG**) for this model is shown in figure 1. The value of the hyperparameters we use are

$$\alpha = 1, \quad (8)$$

$$\xi = 0.0, \quad (9)$$

$$\lambda = 1.0, \quad (10)$$

$$a = 2.0, \quad (11)$$

$$b = 2.0. \quad (12)$$

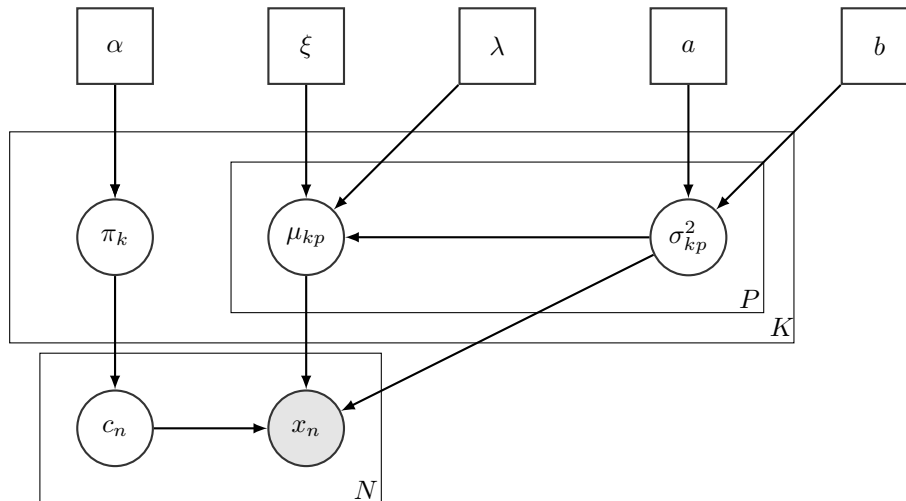


Figure 1: Directed acyclic graph for a mixture of Gaussians with independent features, as used in the simulation study.

2.2 Integrative clustering

We are interested in the use of consensus clustering for integrative methods. We used Multiple Dataset Integration (**MDI**, Kirk et al., 2012) as an example of a Bayesian integrative clustering method. MDI models dataset specific clusterings, in contrast to, for example, Clusternomics (Gabasova et al., 2017) in which a *global clustering* is inferred.

The defining aspect of MDI is the prior on the allocation of the n^{th} item across the L datasets

$$p(c_{n1}, \dots, c_{nL}) \propto \prod_{l=1}^L \pi_{c_{nl}l} \prod_{l=1}^{L-1} \prod_{m=l+1}^L (1 + \phi_{lm} \mathbb{I}(c_{nl} = c_{nm})) \text{ for } n = 1, \dots, N. \quad (13)$$

ϕ_{lm} is the parameter defined by the similarity of the clusterings for the l^{th} and m^{th} datasets and is also sampled in each iteration. As ϕ_{lm} increases more mass is placed on the common partition for these datasets. Conversely, in the limit $\phi_{lm} \rightarrow 0$ we have independent mixture models. In other words, MDI allows datasets with similar clustering of the items to inform the clustering in each other more strongly than the clustering for an unrelated dataset. The DAG for this model for three datasets is shown in figure 2.

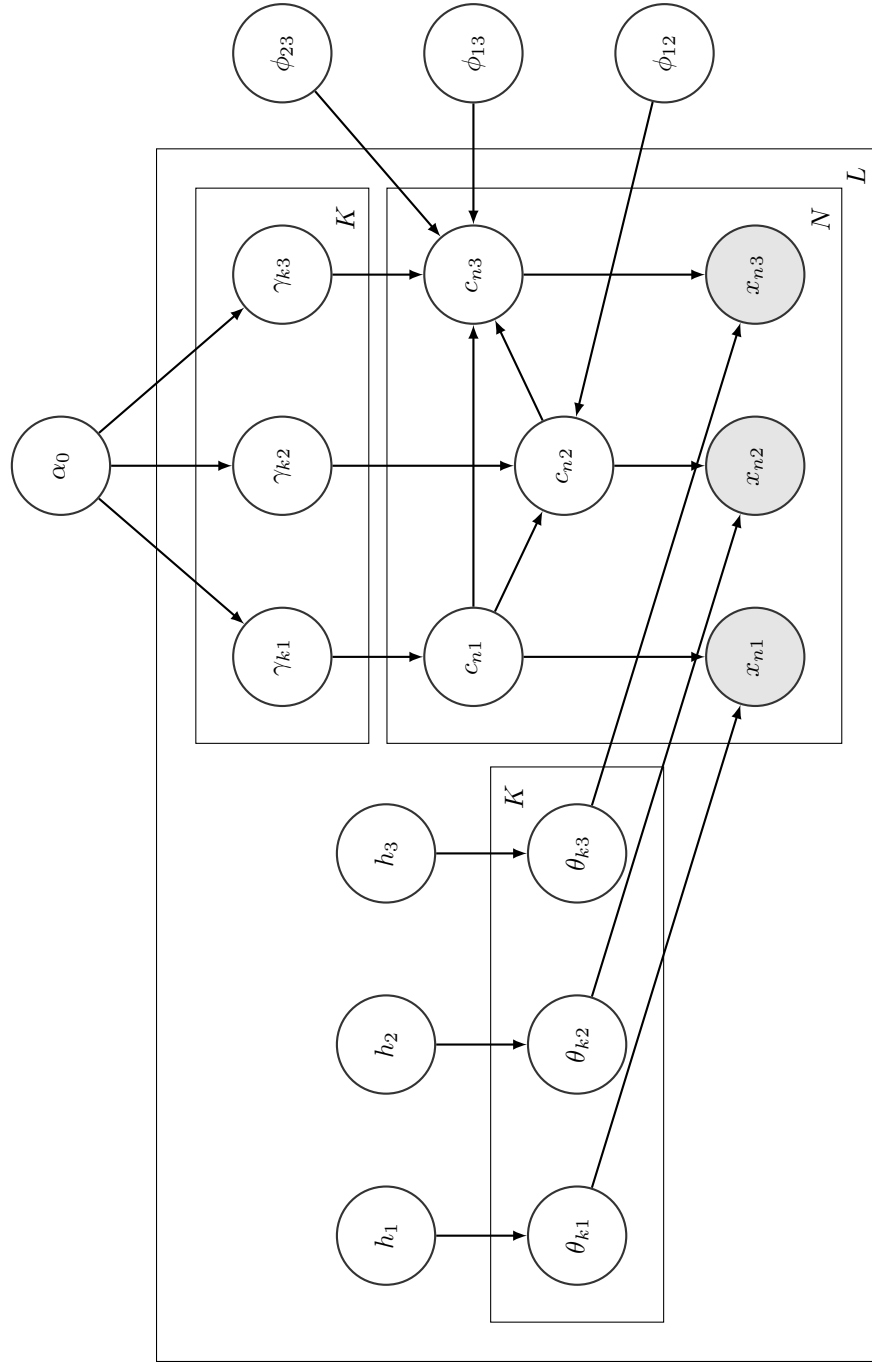


Figure 2: Directed acyclic graph for the Multiple Dataset Integration model for $L = 3$ datasets. h_l is the choice of hyperpriors for the l^{th} dataset.

3 Consensus clustering

Consensus clustering as described by Monti et al. (2003) applies W independent runs of the underlying clustering algorithm to perturbed versions of the dataset and combines the W final partitions in a *consensus matrix* which can be used to infer a final clustering. An outline of this is described in algorithm 1.

The consensus matrix is a symmetric matrix with the $(i, j)^{th}$ entry being the proportions of model runs for which the i^{th} and j^{th} items are clustered together.

Data: $X = (x_1, \dots, x_N)$
Input: A resampling scheme *Resample*
 A clustering algorithm *Cluster*
 Number of resampling iterations W
 Set of cluster numbers to try $\mathcal{K} = \{K_1, \dots, K_{max}\}$
Output: A predicted clustering, \hat{Y}
 The predicted number of clusters present \hat{K}

```

begin
  for  $K \in \mathcal{K}$  do
    /* initialise an empty consensus matrix */
     $\mathbf{M}^{(K)} \leftarrow \mathbf{0}_{N \times N}$ ;
    for  $w = 1$  to  $W$  do
       $X^{(s)} \leftarrow Resample(X)$ ;
      /* Cluster the peturbed dataset, represented in a
         coclustering matrix */
       $\mathbf{B}^{(w)} \leftarrow Cluster(X^{(w)}, K)$ ;
       $\mathbf{M}^{(K)} \leftarrow \mathbf{M}^{(K)} + \mathbf{B}^{(s)}$ ;
    end
     $\mathbf{M}^{(K)} \leftarrow \frac{1}{W} \mathbf{M}^{(K)}$ ;
  end
   $\hat{K} \leftarrow$  best  $K \in \mathcal{K}$  based upon all  $\mathbf{M}^{(K)}$ ;
   $\hat{Y} \leftarrow$  partition  $X$  based upon  $\mathbf{M}^{(\hat{K})}$ ;
end

```

Algorithm 1: Consensus clustering algorithm

To partition X based upon the consensus matrix, we use the R function `maxpear` (Fritsch, 2012). `maxpear` uses a sample average clustering, inferring this by maximising the quantity

$$\frac{\sum_{i < j} \mathbb{I}(c_i^* = c_j^*) p_{ij} - \sum_{i < j} \mathbb{I}(c_i^* = c_j^*) \sum_{i < j} p_{ij} / \binom{N}{2}}{\frac{1}{2} \left[\sum_{i < j} \mathbb{I}(c_i^* = c_j^*) + \sum_{i < j} p_{ij} \right] - \sum_{i < j} \mathbb{I}(c_i^* = c_j^*) \sum_{i < j} p_{ij} / \binom{N}{2}} \quad (14)$$

where p_{ij} is the $(i, j)^{th}$ entry of the consensus matrix (Fritsch et al., 2009).

4 Simulated data

4.1 Scenario description

We defined 12 scenarios to simulate data within to test consensus clustering and to compare it to some alternative tools. Table 1 describes the parameters defining these scenarios and algorithm 2 describes how individual simulations were generated.

Scenario	N	P_s	P_n	K	$\Delta\mu$	σ^2	π
2D	100	2	0	5	3.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
No structure	100	0	2	1	0.0	1	1
Base Case	200	20	0	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Large standard deviation	200	20	0	5	1.0	9	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Large standard deviation	200	20	0	5	1.0	25	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Irrelevant features	200	20	10	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Irrelevant features	200	20	20	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Irrelevant features	200	20	100	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Varying proportions	200	20	0	5	1.0	1	$(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{16})$
Varying proportions	200	20	0	5	0.4	1	$(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{16})$
Small N , large P	50	500	0	5	1.0	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$
Small N , large P	50	500	0	5	0.2	1	$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$

Table 1: Parameters defining the simulation scenarios as used in generating data and labels.

We intend the scenarios to test different aspects of real data or to benchmark performance for comparison in the more challenging situations.

- *2D*: a low dimensional scenario within which we expected **Mclust** to perform well and the long chains to converge and explore the full support of the posterior distribution.
- *No structure*: we included this scenario to reassure fears that consensus clustering has a predilection to finding clusters where none exist (Şenbabaoğlu et al., 2014a,b).
- *Base case*: highly informative datasets within which we expected methods to find the true generating labels quite easily. We included this scenario to benchmark the others that are variations of this setting.
- *Large standard deviation*: these two scenarios investigated the degree of distinction required between clusters for the methods to uncover their structure.
- *Irrelevant features*: we included these scenarios to investigate how robust the methods are to irrelevant features.

Algorithm: Simulation generation

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 ← 0N×P;
/* create a matrix of K means */
μ ← (Δμ, ..., KΔμ);
/* generate the allocation vector */
c ← Sample(1 : K, N, π);
M ← 0N×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 φp = 1 then
    ν ← Permute(μ);
    for n = 1 to N do
      | X(n, p) ← Gaussian(νcn, σ2)
    end
  end
  if φp = 0 then
    for n = 1 to N do
      | X(n, p) ← Gaussian(0, σ2)
    end
  end
end
/* Mean centre and scale the data */
X ← Normalise(X)
```

end

Algorithm 2: 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.

- *Varying proportions*: these scenarios investigated how well each method uncovers clusters when the clusters have significantly different membership counts.
- *Small N, large P*: an investigation of behaviour when the number of features is far greater than the number of items.

4.2 Mclust

We called `Mclust` using the default settings and a range of inputs for the choice of K . We used $K = (2, \dots, \min(\frac{N}{2}, 50))$ to mirror the choice of $K_{max} = 50$ used for the overfitted mixture models (the default in the software we used), with the bound of $\frac{N}{2}$ to avoid fitting 50 clusters in the *Small N, large P* scenario where $N = 50 = K_{max}$. In the *No structure* scenarios we extended to range to $K = (1 \dots, 50)$ to include the correct structure as an option. The model choice was performed using the Bayesian Information Criterion (Schwarz et al., 1978, as implemented in `Mclust`). `Mclust` tries different covariance matrices and thus the model choice is not just between different values of K .

4.3 Bayesian analysis

We use the implementation of Bayesian mixture models in C++ provided by Mason et al. (2016). Rather than directly using a Dirichlet process (Ferguson, 1973) to infer the number of clusters or a mixture that grows and shrinks (Richardson and Green, 1997), this implementation follows the logic of Rousseau and Mengersen (2011) and Van Havre et al. (2015) using an overfitted mixture model to approximate a Dirichlet process. In overfitted mixture models, the number of components, K_{max} , included in the model is set to number far larger than the true number of clusters, K , present.

For each simulation we ran 10 chains for 1 million iterations, keeping every thousandth sample. We discarded the first 10,000 iterations to account for burn-in bias, leaving 990 samples per chain. To check if the chains were converged we used

- the Geweke convergence diagnostic (Geweke et al., 1991) to investigate within-chain stationarity, and
- the potential scale reduction factor (\hat{R} , Gelman et al., 1992) and the Vats-Knudson extension (*stable* \hat{R} , Vats and Knudson, 2018) to check across-chain convergence.

The Geweke convergence diagnostic is a standard Z-score; it compares the sample mean of two sets of samples (in this case buckets of samples from the first half of the samples to the sample mean of the entire second half of samples). It is calculated under the assumption that the two parts of the chain are asymptotically independent and if this assumption holds (i.e. the chain is sampling the same distribution in both samples) than the scores are expected to be standard normally distributed. If a chain's Geweke convergence diagnostic passed

a Shapiro-Wilks test for normality (Shapiro and Wilk, 1965) (based upon a threshold of 0.05), we considered it to have achieved stationarity and included it in the model performance analysis.

\hat{R} is expected to approach 1.0 if the set of chains are converged. Low \hat{R} is not sufficient in itself to claim chain convergence, but values above 1.1 are clear evidence for a lack of convergence (Gelman et al., 2013). Vats and Knudson (2018) show that this threshold is significantly too high (1.01 being a better choice) and propose extensions to \hat{R} that enable a more formal rule for a threshold. We use their method as implemented in the R package `stableGR` (Knudson and Vats, 2020) as the final check of convergence. An example of the \hat{R} series across the 100 simulations for a scenario where chains are well-behaved is shown in figure 3.

We focused upon stationarity of the continuous variables as assessing convergence of the allocation labels is difficult due to label-switching. In our simulations the only recorded continuous variable is the concentration parameter of the Dirichlet distribution for the component weights.

We pooled the samples from the stationary chains and used these to form a PSM. This and the point estimate clustering found by applying the R function `maxpear`. In Bayesian inference, `maxpear` attempts to find the clustering that maximises the Adjusted Rand Index to the true clustering by using an approximation of the expected clustering under the posterior, $\mathbb{E}(c|X)$, believing that this converges to the true clustering. A sample average clustering is used to approximate the expected clustering. This is estimated from the PSM by maximising

$$\frac{\sum_{i < j} \mathbb{I}(c_i^* = c_j^*) p_{ij} - \sum_{i < j} \mathbb{I}(c_i^* = c_j^*) \sum_{i < j} p_{ij} / \binom{N}{2}}{\frac{1}{2} \left[\sum_{i < j} \mathbb{I}(c_i^* = c_j^*) + \sum_{i < j} p_{ij} \right] - \sum_{i < j} \mathbb{I}(c_i^* = c_j^*) \sum_{i < j} p_{ij} / \binom{N}{2}} \quad (15)$$

where p_{ij} is the $(i, j)^{th}$ entry of the PSM (Fritsch et al., 2009). When the chain has converged this maximises the posterior expected ARI to the true clustering.

There are three possibilities to consider the decision to pool the samples across chains under:

- The chains are converged and agree upon the distribution sampled (see figure 4 for an example).
- The chains are not in agreement upon the partition sampled, becoming trapped in different modes. However, a mode does dominate being the mode present in a majority of chains (see figure 5 for an example of this behaviour).
- The chains are not in agreement and no one mode dominates among chains (see figure 6 for an example of this behaviour).

In the first case pooling has no effect upon the predicted clustering compared to using any one chain. In the second case it feels natural that one would use the mode that dominates. Pooling the samples effectively does this for the

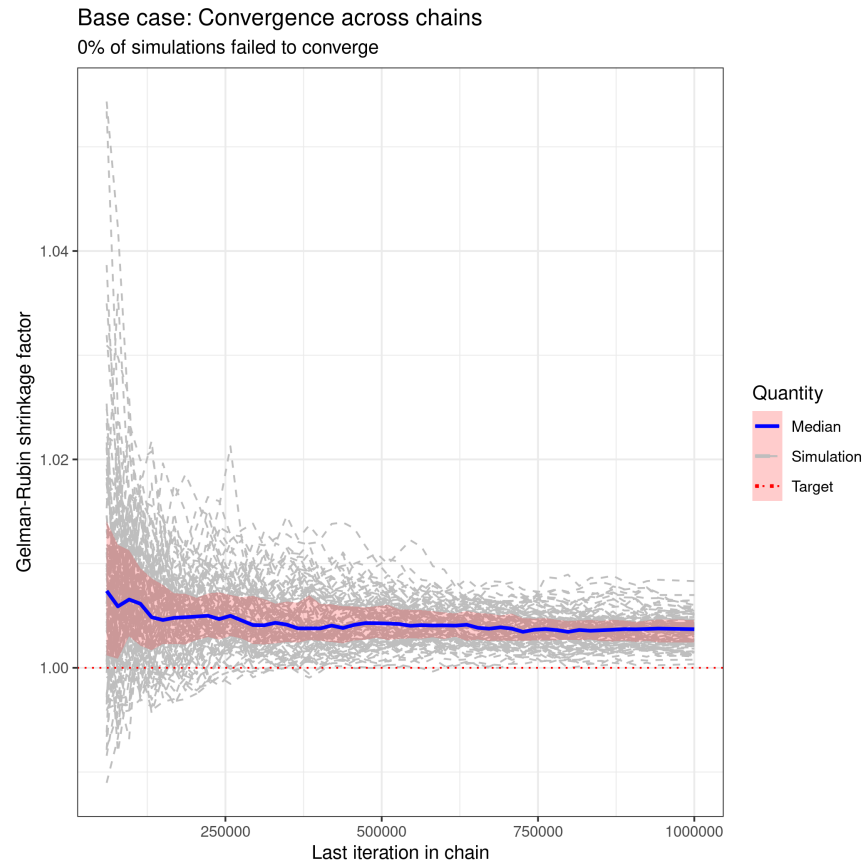


Figure 3: The \hat{R} values for each simulation (in dotted grey), the median value and the interquartile range across simulations. One can see that \hat{R} approaches 1.0, being below 1.01 for every simulation by the end of the chains. The “0% of simulations failed to converge” is a statement based upon the percentage of simulations which passed the test of stable \hat{R} .

Large standard deviation ($\sigma^2 = 9$)

Posterior similarity matrices (simulation 1)

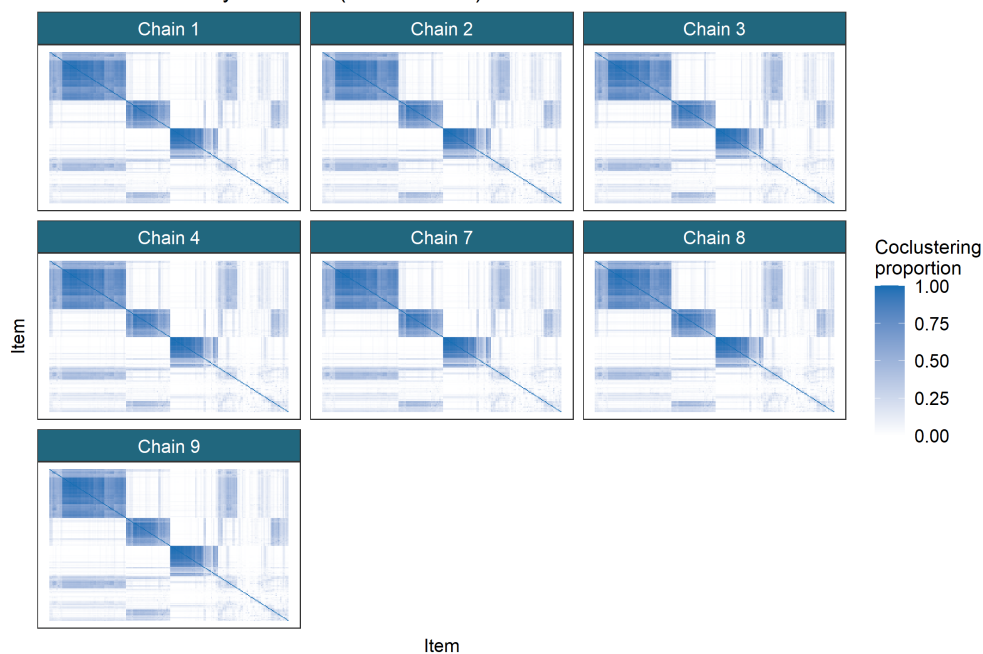


Figure 4: Posterior similarity matrices for the simulation generated using a random seed set to 1 for the first large standard deviation scenario from table 1. This is an example of all stationary chains agreeing in a simulation (and thus pooling of samples is no different to using any choice of chain for the performance analysis). Ordering of rows and columns is defined by hierarchical clustering of the first matrix in the series, in this case that from Chain 1.

predictive performance of the method as the mode with the greatest number of samples across the chains dominates; however, the uncertainty for this mode is increased. In the third case the analysis is non-trivial and further thought, chains and samples would be required. In our simulations this case only arises in the most pathological form in the second *Large N, small P* scenario, where each chain remains trapped in the initial partition. The clustering inferred from any chain is not meaningful being a random clustering; thus the clustering predicted by pooling the PSMs is no more or less relevant as it too is random.

Small N large P ($\Delta\mu = 1.0$)

Posterior similarity matrices (simulation 1)

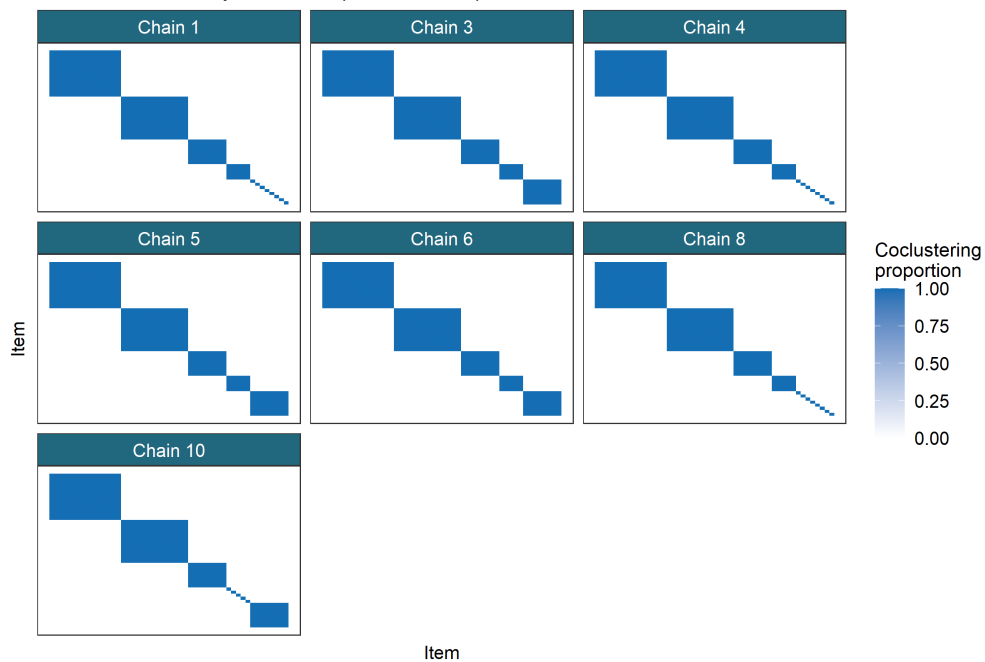


Figure 5: Posterior similarity matrices for the simulation generated using a random seed set to 1 for the first small N , large P scenario from table 1. This is an example of different chains becoming trapped in different modes, but one mode (which does represent the generating structure well) is dominant, being fully present in 3 of the 6 chains, with the two other modes present having significant overlap. Ordering of rows and columns is defined by hierarchical clustering of the first matrix in the series, in this case that from Chain 1.

Small N large P ($\Delta\mu = 0.2$)

Posterior similarity matrices (simulation 1)

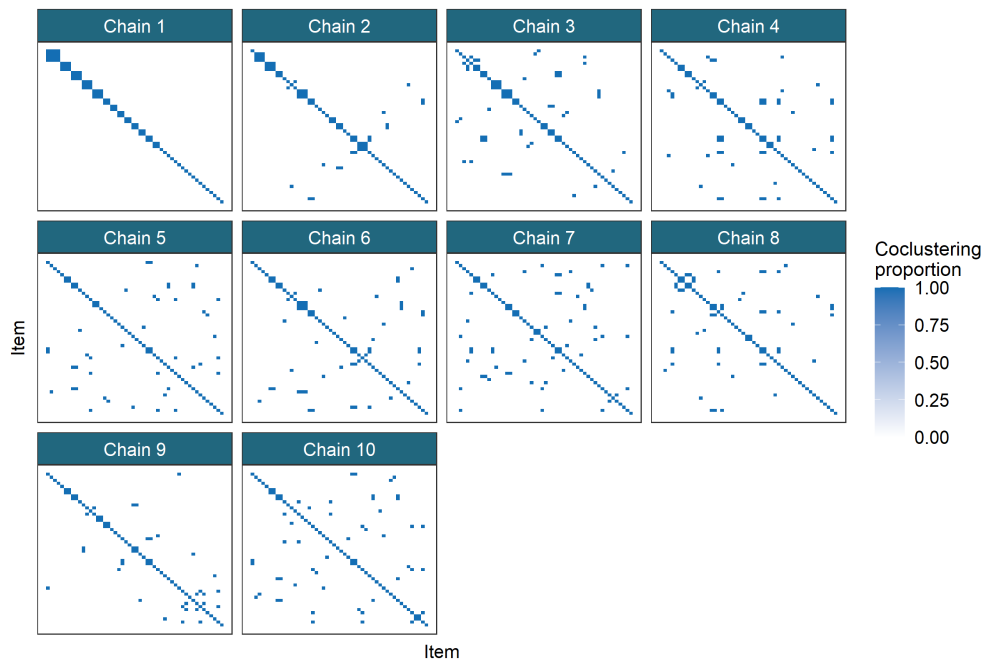


Figure 6: Posterior similarity matrices for the simulation generated using a random seed set to 1 for the second small N , large P scenario from table 1. This is an example of different chains becoming trapped in different modes with no mode being dominant. In this scenario each chain remains trapped in initialisation. Ordering of rows and columns is defined by hierarchical clustering of the first matrix in the series, in this case that from Chain 1.

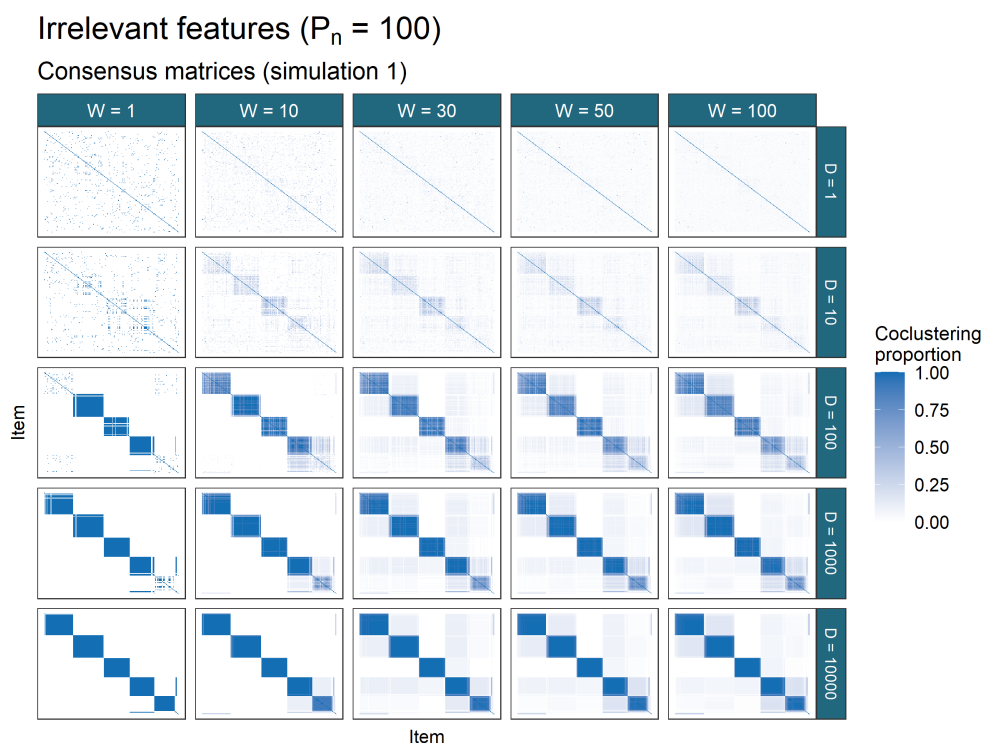


Figure 7: Consensus matrices for the simulation generated using a random seed set to 1 for the third irrelevant features scenario from table 1. D is the individual chain length and W is the number of chains used. In this example there are several modes present (as seen in the entries with values between 0 and 1) but one mode is clearly dominant (the 5 dark squares along the diagonal which correspond closely to the generating labels).

4.4 Consensus clustering analysis

We investigated a range of ensembles, using all combinations of chain depth, $D = \{1, 10, 100, 1000, 10000\}$, and the number of chains, $W = \{1, 10, 30, 50, 100\}$. This gave a total of 25 different ensembles. A consensus matrix was constructed from the samples generated by each ensemble by finding the proportion of samples within which any pair of items are coclustered. An example of the Consensus matrices for each ensemble in a given simulation is shown in figure 7.

4.5 Model performance

The different models (Bayesian (pooled), `Mclust` and the 25 consensus clustering ensembles) were compared under their ability to uncover the generating clustering.

In figure 11 the ARI between the generating labels and the point estimate clustering from each method is shown. For two partitions c_1, c_2 ,

- $ARI(c_1, c_2) = 1.0$: a perfect match between the two partitions,
- $ARI(c_1, c_2) = 0.0$: c_1 is no more similar to c_2 than is expected for a random partition of the data.

In several scenarios `Mclust` performs the best under this metric (e.g. in the scenarios *2D*, *Small N*, *large P* ($\Delta\mu = 0.2$)). However when the number of irrelevant features is large `Mclust` performs less well (see *Irrelevant features* ($P_n = 20$) and ($P_n = 100$)) than the other methods. In the scenario that $P_n = 100$ failing to find structure is not inherently wrong as a majority of the features suggest that there are no subpopulations.

For the ensembles there are two parameters changing between each model, the iteration used to provide the clustering in the ensemble, D , and the number of chains (and hence samples) used, W . In many of the scenarios we find that the benefit of increasing D stabilises by approximately $D = 10$. We believe that in a low-dimensional dataset (such as *2D*), or a highly informative dataset (such as *Base case* or any of the higher dimensional scenarios with no irrelevant features where $\frac{\Delta\mu}{\sigma^2} \geq 1$) the chains quickly find a “sensible” partition of the data and thus increasing the depth within the chain does not increase the probability that any partition sampled will be closer to the generating partition. For example in figure 11 in the *Small N*, *large P* case, the distribution of the ARI across the ensembles for which $D \geq 10$ and $W = 1$ is nearly identical; this suggests that the chain is sampling a very similar partition again and again for 9,990 iterations (and possibly beyond based upon the PSMs shown in figure 5) and it is through adding more chains rather than using particularly long chains that we improve the ability to uncover the generating structure.

We also notice that even if the behaviour has not stabilised for D that the ensemble can uncover meaningful structure. The ARI for the ensembles of short chains can be quite high (as is the case in many of the scenarios). The behaviour of the consensus matrices also shows that low D is not a disqualifier from meaningful inference even if longer chains would be ideal, a result that might be useful in real applications with large datasets and complex models. Consider the consensus matrices in figure 7, it can be seen that the behaviour has not stabilised before $D = 10000$ (and possibly there is still some benefit in increasing D beyond this value), but the structure being uncovered when there is a sufficient number of chains and D is small does correspond to the structure uncovered in the largest and deepest ensemble. We believe that the order in which components merge and items are co-clustered varies depending on initialisation, and thus if the chain is not sufficiently deep that all of the final mergings have occurred that a sufficiently large ensemble can still perform meaningful inference of the subpopulation structure despite the poor performance of any individual model. Even though each learner probably has too many clusters for small D the consensus among them will have less if the individual learners have low correlation between their partitions (something we might expect if the chains are

stopped very early). This is why the entries of the consensus matrix for $D = 100$ and $W = 100$ in figure 7 are more pale than in deeper ensembles; very few items correctly (possibly none) cocluster in every partition, it is only in observing the consensus that the global structure of interest emerges. Thus if there is some limit to the length of chains available for an analysis (e.g. computational or temporal constraints) than the inference obtained from the shorter chains can still be meaningful, with the caveat that the point clustering might have more clusters than the same analysis with longer chains would provide. Additional post-hoc merging of some clusters might be necessary in this case.

In contrast, when the dataset is sparse or contains many irrelevant features, we believe that deeper chains are required to reach this steady-state sampling where no single sample is expected to be better than any other (see the *Irrelevant features* ($P_n = 100$) facet of figure 11).

In some scenarios no method is successful in uncovering the generating labels. In the *Large standard deviation* ($\sigma^2 = 25$) and *Small N , large P* ($\Delta\mu = 0.2$) this is due to the lack of signal - the clusters overlap so significantly that it is not possible for any of these methods to uncover much of the generating structure. In the *No structure* case it is different (although `Mclust` does perform well here). In this case all items are generated from a common distributions. For the Bayesian chains and the ensembles, a clustering of singletons is predicted; each item is allocated a unique label (see figures 8 and 9). While failing to perform well under the ARI, this is a sensible result. Rather than indicating (as we did with the shared label) that no item is particularly distinct from the others and thus all share a common label, this clustering of singletons states that no item is more similar to any other and thus no two items should cluster together. It is an alternative statement of the same result, i.e. that there is no evidence for subpopulation structure. We consider this evidence that an ensemble of Bayesian mixture models is not as susceptible to predicting labels than an ensemble based upon K -means clustering as in Şenbabaoğlu et al. (2014a,b).

Increasing W is also required when the dimensionality of the dataset is large. In this case it is due to individual chains exploring only a single mode (as can be seen in figure 5 where each chain appears to sample only a single partition). In this example where each sample is a partition that appears to be a mode in the posterior distribution of the allocation vector from very early in the chain (based upon the stable performance for $D \geq 10$), increasing W allows each chain to “vote” on which mode is the global mode, as we believe that the mode that attracts the most chains is the global mode (although in real datasets the number of chains required might be greater than in our simulations). An example of this behaviour may be seen in figure 10.

In figure 11, limiting behaviour for increases of W and D can be seen for the ensemble. For most simulations there is no change in performance for greater choices of W and D after some stabilising values.

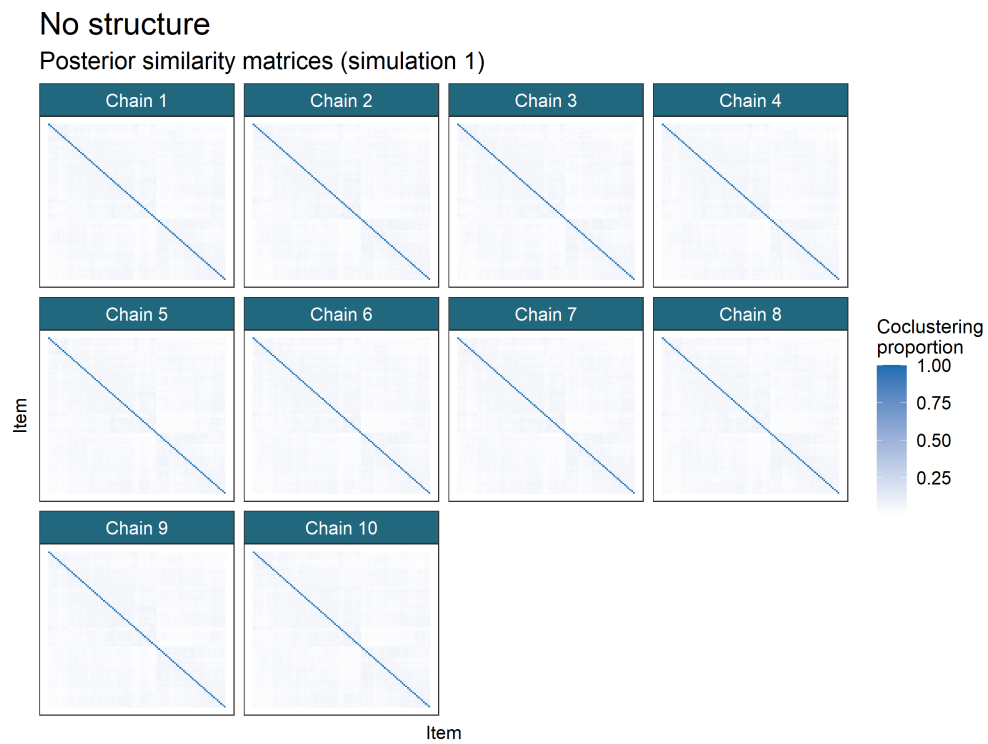


Figure 8: Posterior similarity matrices for simulation 1 of the *No structure* scenario. Each item is allocated to a singleton.

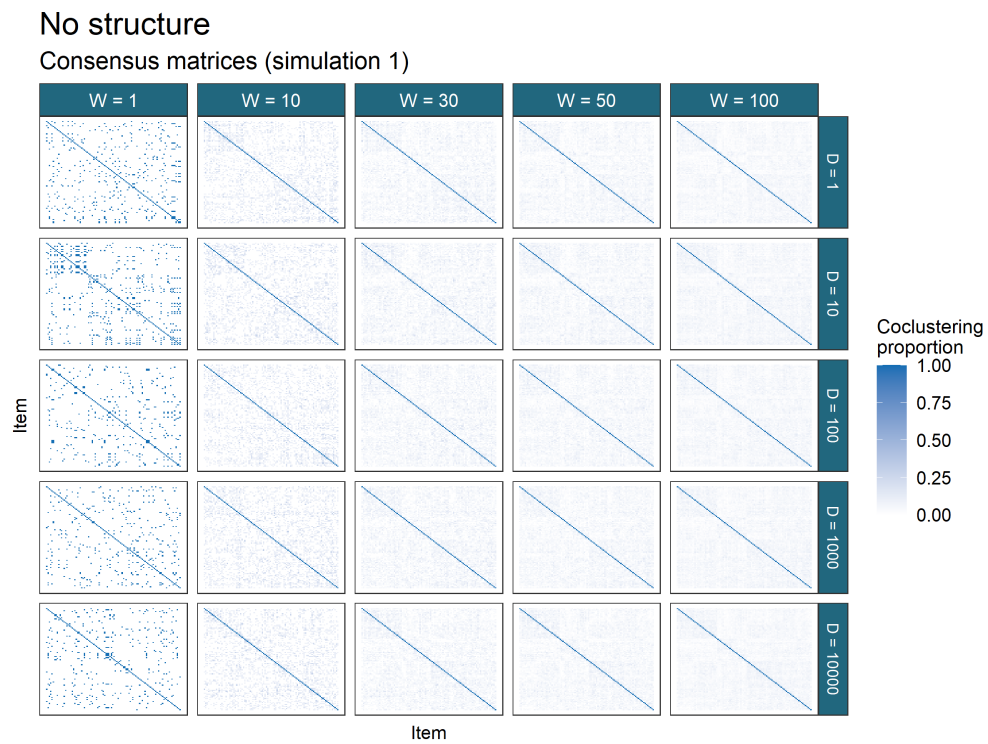


Figure 9: Consensus matrices for simulation 1 of the *No structure* scenario. Each item is allocated to a singleton in many of the Consensus matrices.

Small N large P ($\Delta\mu = 1.0$)

Consensus matrices (simulation 1)

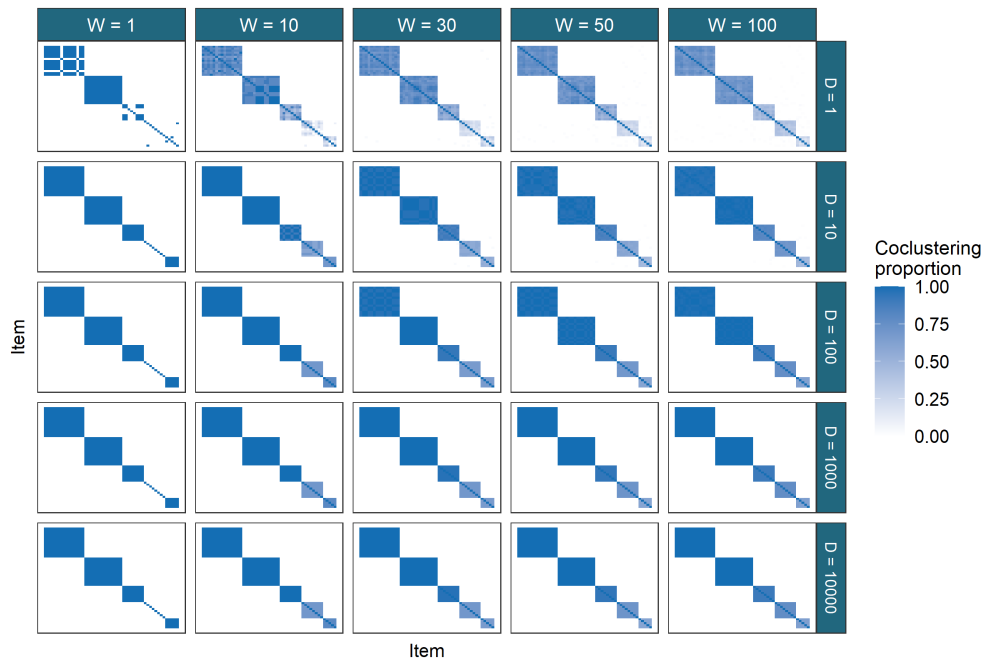


Figure 10: Consensus matrices for simulation 1 of the first *Large N, small P* scenario. One can see that by iteration ten the sample being drawn is from the mode (for $W = 1, D = 10$), and that an ensemble of chains does find structure that recalls the generating labels (see figure 11, the ARI for $CC(10, s)$ is 1.0 for $s > 1$, meaning that the true labels perfectly align with those predicted by the consensus matrix).

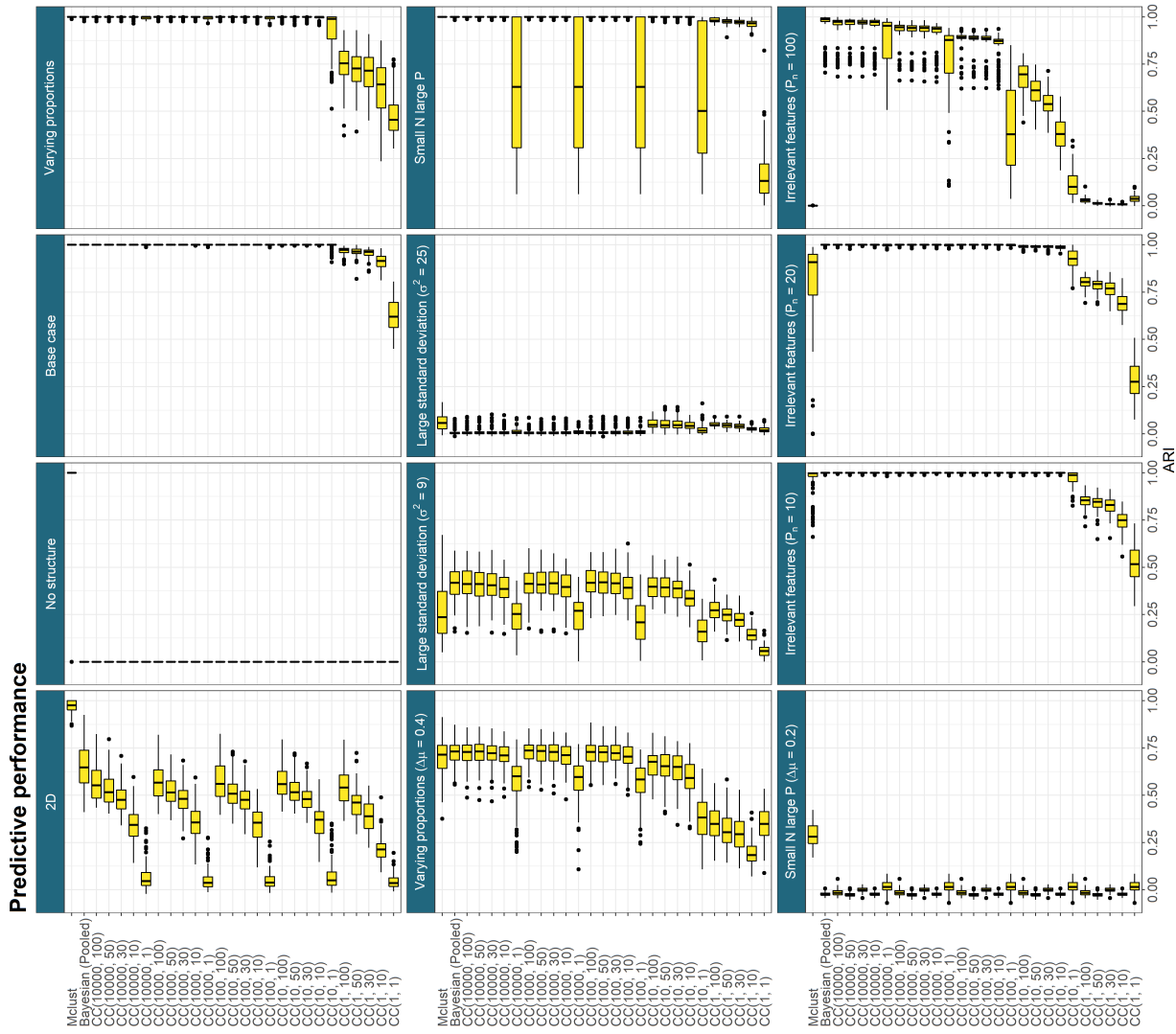


Figure 11: Predictive performance across all simulations. $CC(D, W)$ denotes consensus clustering using the D^{th} sample from W different chains. In the cases where the generating structure is not exactly found, increasing D and W sees some improvement in the ARI between the truth and the predicted clusterings before some limiting behaviour emerges and and further increase appears to have no change in the performance.

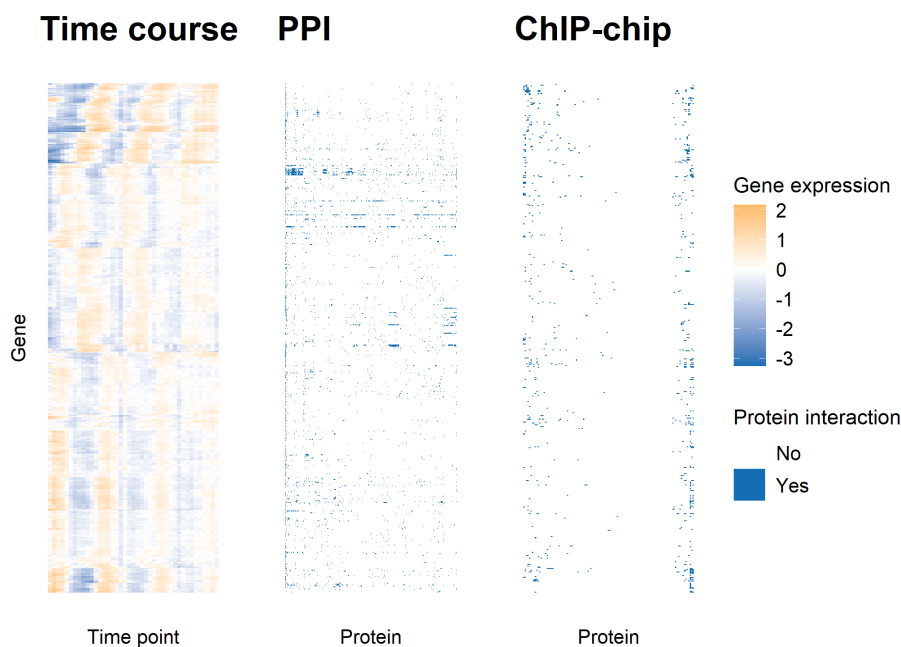


Figure 12: Heatmap of the yeast datasets. Each plot has a common row order corresponding to the gene products being clustered. This order was decided by a hierarchical clustering of the rows of the time course expression matrix. The time course data is associated with the “Gene expression” legend and the ChIP-chip and PPI data with “Protein interaction” legend.

5 Multi-omics analysis of the cell cycle in budding yeast

We chose our three datasets (shown in figure 12) to perform an integrative analysis as many of the protein encoding genes in the mitotic cell cycle have well studied genomic binding sites with mapped transcription factors (TFs) that control phase-specific expression (Cho et al., 1998; Spellman et al., 1998); thus the inclusion of the ChIP-chip data means that the clusters that align across the datasets should include well studied regulatory proteins and thus be of biological interest. If a cluster of genes are similarly expressed in the time course, share associated regulatory protein in the ChIP-chip and are associated with common protein complexes in the PPI data, then this implies a gene set with strong biological significance.

In contrast, if we cluster the time course dataset alone, any clusters that we find are defined by correlation across time. This might be assumed to be driven by shared regulatory mechanisms, but other sources of structure might

be encouraging this, even experimental error. However, if a cluster aligns across both the time course dataset and the ChIP-chip dataset we can be more certain that these genes are part of some regulatory network; if this cluster also emerges in the PPI dataset we might believe that the genes are co-regulated as part of the formation of some protein complex. Furthermore, this integrative aspect means that clusters that might merge in the time course dataset due to similar periodicity in a standalone analysis might remain separate due to different associated transcription factors in the ChIP-chip dataset.

Thus we performed an integrative analysis using MDI to avoid aggressive assumptions about either the biology defining any clusters and modelling assumptions about the latent structure.

We expect that the complexity of this data and model means that the time required for convergence of the MCMC algorithm would be very large. We avoid this problem by using consensus clustering of MDI, instead basing our final ensemble choice on the stopping rule described in the main paper.

The datasets were modelled using a mixture of Gaussian processes in the time course dataset and Multinomial distributions in the ChIP-chip and PPI datasets. To ensure that our mixture model is initially overfitted we set $K_{max} = 275 \approx \frac{N}{2}$, and following from this the point estimate was inferred from the consensus matrix using `maxpear` as in the simulated data except we set `k.max` = 275.

5.1 Consensus clustering analysis

We include the consensus matrices for each dataset for a range of ensembles for further evidence that the ensemble was stable for the 10,000th iteration from 1,000 chains in figures 13, 14 and 15. In these figures, there is no strong change between the consensus matrices for $D = 5001$ and $D = 10001$.

We wish to identify groups of genes that tend to be grouped together in multiple datasets. We focus upon the genes that tend to have the same cluster label in multiple datasets, those which have a common label across some set of datasets in more than half of the observed clusterings, or $\hat{P}(c_{nl} = c_{nm}) > 0.5$, where c_{nl} denotes the cluster label of gene n in dataset l . This based upon the the concept of *fused genes* proposed by Savage et al. (2010) and used by Kirk et al. (2012), but to avoid confusion due to other possible ideas of fused genes (e.g. those that contribute to a common protein complex, the behaviour of TFs upon a gene) we avoid this term. These genes with common clustering across datasets are those most affected by the integrative aspect of the analysis and therefore we focus upon these in the our cluster analysis. In our case we have the possible sets of:

- {Time course}, {ChIP-chip}, {PPI},
- {Time course, ChIP-chip}, {Time course, PPI}, {ChIP-chip, PPI}, and
- {Time course, ChIP-chip, PPI}.

The number of genes meeting this criteria between any two datasets is indicative of how strongly they influence each other and is expected to align with

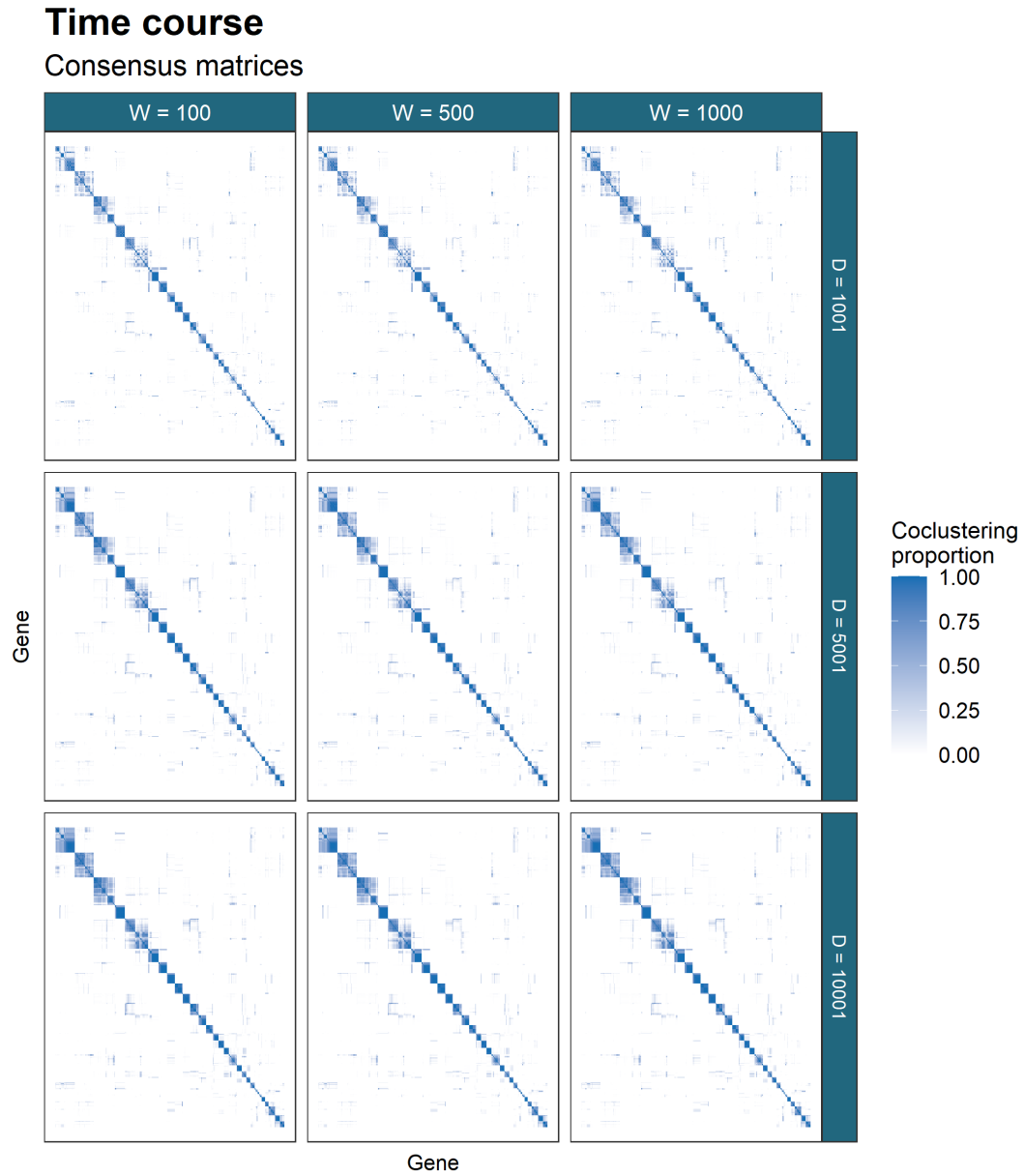


Figure 13: Consensus matrices for different ensembles of MDI for the time course data. This dataset has stable clustering across the different choices of number of chains, W , and chain depth, D , with some components merging as the chain depth increases.

ChIP-chip

Consensus matrices

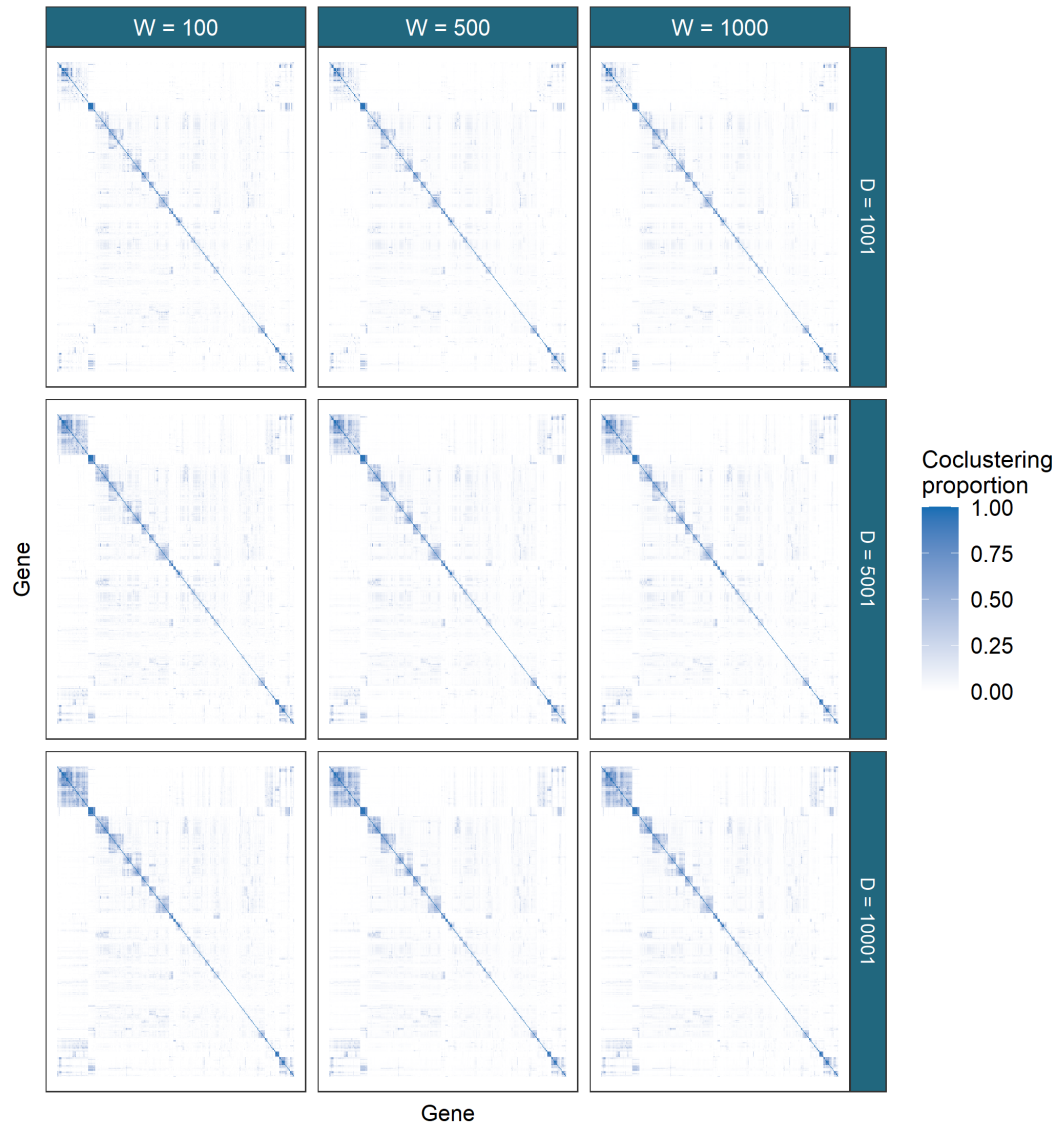


Figure 14: The ChIP-chip dataset is more sparse than the time course data. In keeping with the results from the simulations for mixture models, deeper chains are required for better performance. It is only between $D = 5,001$ and $D = 10,001$ that no change in the clustering can be observed and the result is believed to be stable. In this dataset the number of chains used, W , appears relatively unimportant, with similar results for $W = 100, 500, 1000$.

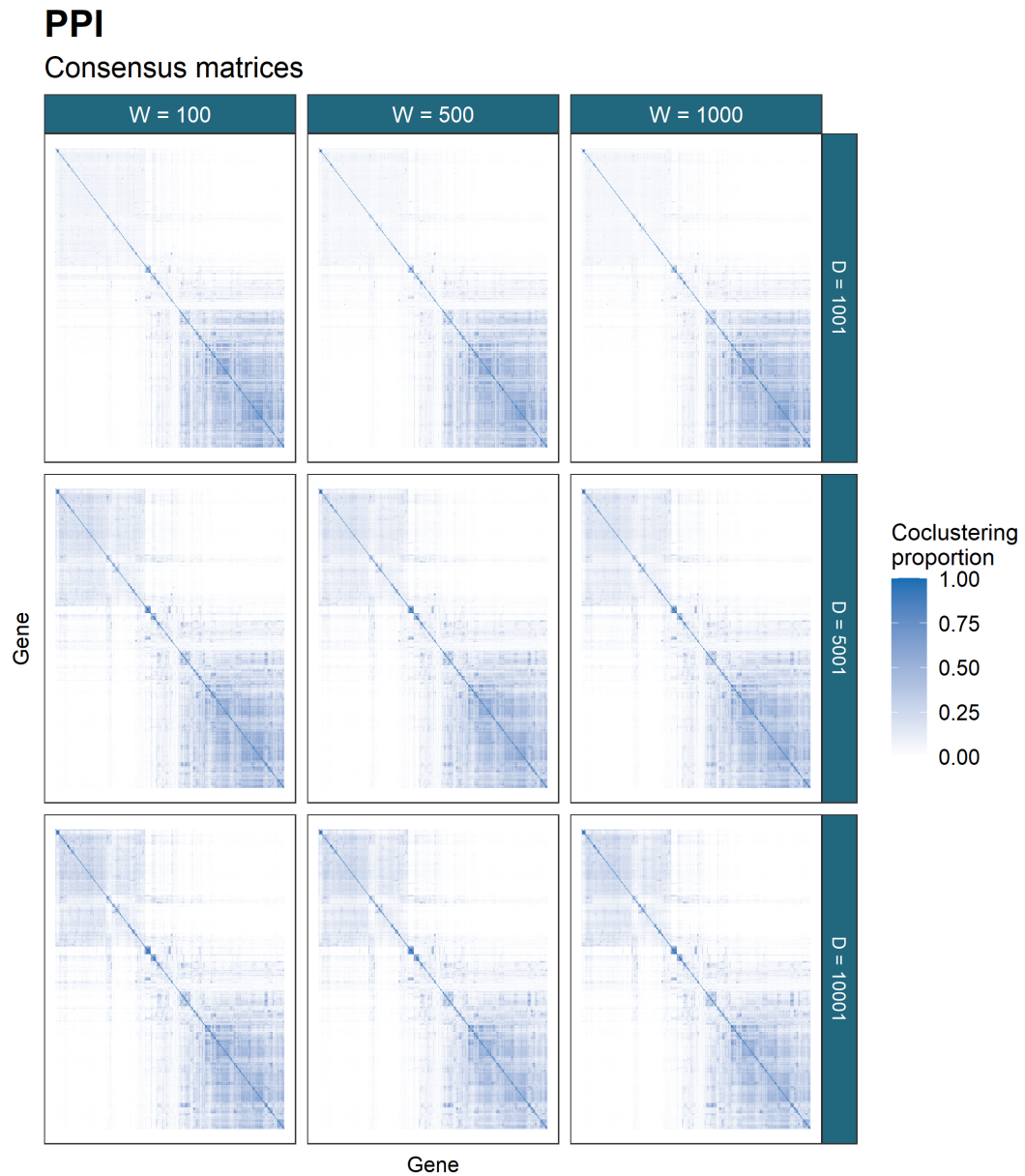


Figure 15: The PPI dataset has awkward characteristics for modelling. A wide, sparse dataset it is chain depth that we found to be the most important parameter for the ensemble. Similar to the results in figure 14, the matrices only stabilise from $D = 5001$ to $D = 10001$.

the ϕ_{lm} parameters from the MDI model. We find the following number of unique genes integrated between each combination of datasets:

- Time course + ChIP-chip + PPI: 56,
- Time course + ChIP-chip: 205 (261 including the 56 integrated across all datasets),
- Time course + PPI: 12 (68),
- ChIP-chip + PPI: 43 (99).

This shows that the time course and ChIP-chip datasets contain very similar structure, the ChIP-chip and PPI datasets have some similarity but significantly less and the time course and PPI datasets have less shared signal again.

Compare this to the original analysis of this data in Kirk et al. (2012), where the number of such genes in each combination is:

- Time course + ChIP-chip + PPI: 16,
- Time course + ChIP-chip: 32 (48),
- Time course + PPI: 16 (32),
- ChIP-chip + PPI: 15 (31).

Our analysis has found significantly more shared structure.

5.1.1 Time course ChIP-chip analysis

We focus upon the dataset pairing of time course + ChIP-chip within the integrative analysis as the combination with the greatest number of genes with shared clustering. We show these genes grouped by their inferred cluster in figure 16. In this plot we exclude the 15 clusters where more than half of the member genes have no interactions in the ChIP-chip data and any clusters of one. We find that a small number of transcription factors dominate, with different combinations emerging across the 10 clusters shown here in table 2. Many of these 10 correspond to transcription factors that are well known to regulate cell cycle expression, namely MBP1, SWI4, SWI6, MCM1, FKH1, FKH2, NDD1, SWI5, and ACE2 (Simon et al., 2001).

Table 2: Table of transcription factors prominent in clusters of genes with shared labels for a majority of samples for the time course and ChIP-chip datasets.

Gene	Name	Description
------	------	-------------

YLR131C	ACE2	Transcription factor required for septum destruction after cytokinesis; phosphorylation by Cbk1p blocks nuclear exit during M/G1 transition; phosphorylation by cyclins Cdc28p and Pho85p prevents nuclear import during cell cycle phases other than cytokinesis; part of RAM network that regulates cellular polarity and morphogenesis; ACE2 has a paralog, SWI5, that arose from the whole genome duplication
YPL049C	DIG1	MAP kinase-responsive inhibitor of the Ste12p transcription factor; involved in the regulation of mating-specific genes and the invasive growth pathway; Dig1p and paralog Dig2p bind to Ste12p
YIL131C	FKH1	Forkhead family transcription factor; evolutionarily conserved lifespan regulator; binds multiple chromosomal elements with distinct specificities, cell cycle dynamics; regulates transcription elongation, chromatin silencing at mating loci, expression of G2/M phase genes
YNL068C	FKH2	Forkhead family transcription factor; rate-limiting activator of replication origins; evolutionarily conserved regulator of lifespan; binds multiple chromosomal elements with distinct specificities, cell cycle dynamics; positively regulates transcriptional elongation; negative role in chromatin silencing at HML and HMR; major role in expression of G2/M phase genes
YDL056W	MBP1	Transcription factor; involved in regulation of cell cycle progression from G1 to S phase, forms a complex with Swi6p that binds to MluI cell cycle box regulatory element in promoters of DNA synthesis genes
YMR043W	MCM1	Transcription factor; involved in cell-type-specific transcription and pheromone response; plays a central role in the formation of both repressor and activator complexes; involved in the transcription of some M/G1 genes Simon et al. (2001).
YOR372C	NDD1	Transcriptional activator essential for nuclear division; essential component of the mechanism that activates the expression of a set of late-S-phase-specific genes; turnover is tightly regulated during cell cycle and in response to DNA damage
YHR084W	STE12	Transcription factor that is activated by a MAPK signaling cascade; activates genes involved in mating or pseudohyphal/invasive growth pathways; cooperates with Tec1p transcription factor to regulate genes specific for invasive growth

YER111C	SWI4	DNA binding component of the SBF complex (Swi4p-Swi6p); a transcriptional activator that in concert with MBF (Mbp1-Swi6p) regulates late G1-specific transcription of targets including cyclins and genes required for DNA synthesis and repair
YDR146C	SWI5	Transcription factor that recruits Mediator and Swi/Snf complexes; activates transcription of genes expressed at the M/G1 phase boundary and in G1 phase; required for expression of the HO gene controlling mating type switching; localization to nucleus occurs during G1 and appears to be regulated by phosphorylation by Cdc28p kinase; SWI5 has a paralog, ACE2, that arose from the whole genome duplication
YLR182W	SWI6	Transcription cofactor; forms complexes with Swi4p (SBF) and Mbp1p (MBF) to regulate transcription at the G1/S transition (Simon et al., 2001); involved in meiotic gene expression; also binds Stb1p to regulate transcription at START; also required for the unfolded protein response, independently of its known transcriptional coactivators
YBR083W	TEC1	Transcription factor targeting filamentation genes and Ty1 expression; Ste12p activation of most filamentation gene promoters depends on Tec1p and Tec1p transcriptional activity is dependent on its association with Ste12p; binds to TCS elements upstream of filamentation genes, which are regulated by Tec1p/Ste12p/Dig1p complex; competes with Dig2p for binding to Ste12p/Dig1p; positive regulator of chronological life span
YML027W	YOX1	Homeobox transcriptional repressor; binds to Mcm1p and to early cell cycle boxes (ECBs) in the promoters of cell cycle-regulated genes expressed in M/G1 phase; phosphorylated by the cyclin Cdc28p; relocalizes from nucleus to cytoplasm upon DNA replication stress

These regulatory proteins are found in different combinations across the clusters. Based upon these combinations we associate each cluster with phases of the cell cycle and or some specific processes.

- Cluster 1: both ACE2 and SWI5 emerge. These regulate specific genes at the end of M and early G1 (McBride et al., 1999; Simon et al., 2001).
- Cluster 2: SWI5. This is similar to cluster 1, as ACE2 is a paralog of SWI5; therefore associated with M/G1. Furthermore, inspection of the expression in the timecourse data shows that the members of cluster 2

largely differentiate from those of cluster 1 based upon amplitude, not periodicity, suggesting that these clusters could be merged.

- Cluster 5: MBP1, SWI4 and SWI6. The SBF complex (Swi4p-Swi6p) is a transcriptional activator that in concert with MBF (Mbp1-Swi6p) regulates late G1-specific transcription of targets including cyclins and genes required for DNA synthesis and repair, controlling the transition to S phase (Simon et al., 2001; Iyer et al., 2001; Aligianni et al., 2009).
- Cluster 9: MBP1 and SWI6. These combine to form MBF, which regulates DNA replication and repair (Iyer et al., 2001).
- Cluster 11: DIG1, SWI4, SWI6, and STE12 emerge in all members with some having associations with TEC1. TEC1 and STE12, controls development, including cell adhesion and filament formation and is negatively regulated by DIG1 and DIG2 (van der Felden et al., 2014).
- Cluster 12: MBP1, SWI4 and SWI6. Similar to cluster 5 in both the time course and ChIP-chip datasets and thus G1/S phase.
- Cluster 16: some MBP1, SWI4 and SWI6. The constituents of this cluster are largely associated with proteins contributing to histones H1, H2A, H2B, H3 and H4, suggesting an S-phase cluster (Ewen, 2000).
- Cluster 17: FKH1 and FKH2. Fkh1p and Fkh2p are required for cell-cycle regulation of transcription during G2/M (Kumar et al., 2000).
- Cluster 20: NDD1 and MCM1 with some FKH2. Mcm1, together with Fkh1 or Fkh2, recruits the Ndd1 protein in late G2, and thus controls the transcription of G2/M genes (Simon et al., 2001; Koranda et al., 2000).
- Cluster 26: YOX1 and MCM1. YOX1 binds to Mcm1p and to early cell cycle boxes (ECBs) in the promoters of cell cycle-regulated genes expressed in M/G1 phase (Pramila et al., 2002).

Gene	Name	Cluster	Description
YJL115W	ASF1	9	Nucleosome assembly factor; involved in chromatin assembly, disassembly; required for buffering mRNA synthesis rate against gene dosage changes in S phase
YLR103C	CDC45	9	DNA replication initiation factor; recruited to MCM pre-RC complexes at replication origins; recruits elongation machinery; binds tightly to ssDNA, which disrupts interaction with the MCM helicase and stalls it during replication stress; mutants in human homolog may cause velocardiofacial and DiGeorge syndromes

YPL241C	CIN2	9	GTPase-activating protein (GAP) for Cin4p; tubulin folding factor C involved in beta-tubulin (Tub2p) folding; mutants display increased chromosome loss and benomyl sensitivity; human homolog RP2 complements yeast null mutant
YPR175W	DPB2	9	Second largest subunit of DNA polymerase II (DNA polymerase epsilon); required for maintenance of fidelity of chromosomal replication; essential motif in C-terminus is required for formation of the four-subunit Pol epsilon; expression peaks at the G1/S phase boundary; Cdc28p substrate
YIL026C	IRR1	9	Subunit of the cohesin complex; which is required for sister chromatid cohesion during mitosis and meiosis and interacts with centromeres and chromosome arms
YCL061C	MRC1	9	S-phase checkpoint protein required for DNA replication; couples DNA helicase and polymerase; defines a novel S-phase checkpoint with Hog1p that coordinates DNA replication and transcription upon osmostress; protects uncapped telomeres; Dia2p-dependent degradation mediates checkpoint recovery; mammalian claspin homolog; subunit of a replication-pausing checkpoint complex, Tof1p-Mrc1p-Csm3p; checkpoint-mediator protein that functions during DNA replication and activates the effector kinase Rad53 (Bando et al., 2009); human ATR homolog (Lao et al., 2018)
YDR097C	MSH6	9	Protein required for mismatch repair in mitosis and meiosis; forms a complex with Msh2p to repair both single-base and insertion-deletion mispairs; also involved in interstrand cross-link repair; potentially phosphorylated by Cdc28p
YNL102W	POL1	9	Catalytic subunit of the DNA polymerase I alpha-primase complex; required for the initiation of DNA replication during mitotic DNA synthesis and premeiotic DNA synthesis
YBL035C	POL12	9	B subunit of DNA polymerase alpha-primase complex; required for initiation of DNA replication during mitotic and premeiotic DNA synthesis; also functions in telomere capping and length regulation
YKL113C	RAD27	9	5' to 3' exonuclease, 5' flap endonuclease; required for Okazaki fragment processing and maturation, for long-patch base-excision repair and large loop repair (LLR), ribonucleotide excision repair
YPL153C	RAD53	9	DNA damage response protein kinase; required for cell-cycle arrest, regulation of copper genes in response to DNA damage; human homolog CHEK2 implicated in breast cancer can complement yeast null mutant

YAR007C	RFA1	9	Subunit of heterotrimeric Replication Protein A (RPA); RPA is a highly conserved single-stranded DNA binding protein involved in DNA replication, repair, and recombination; RPA protects against inappropriate telomere recombination, and upon telomere uncapping, prevents cell proliferation by a checkpoint-independent pathway; role in DNA catenation/decatenation pathway of chromosome disentangling; relocalizes to the cytosol in response to hypoxia
YNL312W	RFA2	9	Subunit of heterotrimeric Replication Protein A (RPA); RPA is a highly conserved single-stranded DNA binding protein involved in DNA replication, repair, and recombination; RPA protects against inappropriate telomere recombination, and upon telomere uncapping, prevents cell proliferation by a checkpoint-independent pathway; in concert with Sgs1p-Top2p-Rmi1p, stimulates DNA catenation/decatenation activity of Top3p; protein abundance increases in response to DNA replication
YAR008W	SEN34	9	Subunit of the tRNA splicing endonuclease; tRNA splicing endonuclease (Sen complex) is composed of Sen2p, Sen15p, Sen34p, and Sen54p; Sen complex also cleaves the CBP1 mRNA at the mitochondrial surface; Sen34p contains the active site for tRNA 3' splice site cleavage and has similarity to Sen2p and to Archaeal tRNA splicing endonuclease
YJL074C	SMC3	9	Subunit of the multiprotein cohesin complex; required for sister chromatid cohesion in mitotic cells; also required, with Rec8p, for cohesion and recombination during meiosis; phylogenetically conserved SMC chromosomal ATPase family member
YNL273W	TOF1	9	Subunit of a replication-pausing checkpoint complex; Tof1p-Mrc1p-Csm3p acts at the stalled replication fork to promote sister chromatid cohesion after DNA damage, facilitating gap repair of damaged DNA; interacts with the MCM helicase; checkpoint-mediator protein that functions during DNA replication and activates the effector kinase RAD53 (Bando et al., 2009); human ATM homolog (Lao et al., 2018)
YMR215W	GAS3	16	Putative 1,3-beta-glucanosyltransferase; has similarity to other GAS family members; low abundance, possibly inactive member of the GAS family of GPI-containing proteins; localizes to the cell wall; mRNA induced during sporulation

YBR009C	HHF1	16	Histone H4; core histone protein required for chromatin assembly and chromosome function; one of two identical histone proteins (see also HHF2); contributes to telomeric silencing; N-terminal domain involved in maintaining genomic integrity
YNL030W	HHF2	16	Histone H4; core histone protein required for chromatin assembly and chromosome function; one of two identical histone proteins (see also HHF1); contributes to telomeric silencing; N-terminal domain involved in maintaining genomic integrity
YPL127C	HHO1	16	Histone H1, linker histone with roles in meiosis and sporulation; decreasing levels early in sporulation may promote meiosis, and increasing levels during sporulation facilitate compaction of spore chromatin; binds to promoters and within genes in mature spores; may be recruited by Ume6p to promoter regions, contributing to transcriptional repression outside of meiosis; suppresses DNA repair involving homologous recombination
YBR010W	HHT1	16	Histone H3; core histone protein required for chromatin assembly, part of heterochromatin-mediated telomeric and HM silencing; one of two identical histone H3 proteins (see HHT2); regulated by acetylation, methylation, and phosphorylation; H3K14 acetylation plays an important role in the unfolding of strongly positioned nucleosomes during repair of UV damage
YNL031C	HHT2	16	Histone H3; core histone protein required for chromatin assembly, part of heterochromatin-mediated telomeric and HM silencing; one of two identical histone H3 proteins (see HHT1); regulated by acetylation, methylation, and phosphorylation; H3K14 acetylation plays an important role in the unfolding of strongly positioned nucleosomes during repair of UV damage
YDR225W	HTA1	16	Histone H2A; core histone protein required for chromatin assembly and chromosome function; one of two nearly identical subtypes (see also HTA2); DNA damage-dependent phosphorylation by Mec1p facilitates DNA repair; acetylated by Nat4p; N-terminally propionylated in vivo
YBL003C	HTA2	16	Histone H2A; core histone protein required for chromatin assembly and chromosome function; one of two nearly identical (see also HTA1) subtypes; DNA damage-dependent phosphorylation by Mec1p facilitates DNA repair; acetylated by Nat4p

YDR224C	HTB1	16	Histone H2B; core histone protein required for chromatin assembly and chromosome function; nearly identical to HTB2; Rad6p-Bre1p-Lge1p mediated ubiquitination regulates reassembly after DNA replication, transcriptional activation, meiotic DSB formation and H3 methylation
YBL002W	HTB2	16	Histone H2B; core histone protein required for chromatin assembly and chromosome function; nearly identical to HTB1; Rad6p-Bre1p-Lge1p mediated ubiquitination regulates reassembly after DNA replication, transcriptional activation, meiotic DSB formation and H3 methylation
YNR009W	NRM1	16	Transcriptional co-repressor of MBF-regulated gene expression; Nrm1p associates stably with promoters via MCB binding factor (MBF) to repress transcription upon exit from G1 phase
YDR113C	PDS1	16	Securin; inhibits anaphase by binding separin Esp1p; blocks cyclin destruction and mitotic exit, essential for meiotic progression and mitotic cell cycle arrest; localization is cell-cycle dependent and regulated by Cdc28p phosphorylation

Table 3: Description of the genes with common labelling across the time course and ChIP-chip datasets from clusters 9 and 16.

5.2 Bayesian analysis

We wished to compare our results from consensus clustering to a conventional Bayesian approach. We ran 10 chains of MDI for 36 hours saving every thousandth sample. This resulted in chains of varying length. We reduced the chains to 666 samples as this was the number of samples achieved by the shortest chain. Similar to section 4.3 these chains were then investigated for

- within-chain stationarity using the Geweke convergence diagnostic (Geweke et al., 1991), and
- across-chain convergence using \hat{R} (Gelman et al., 1992) and the Vats-Knudson extension (*stable* \hat{R} , Vats and Knudson, 2018).

Again we focus upon stationarity of the continuous variables. In the implementation of MDI we used, the recorded continuous variables are the concentration parameters of the Dirichlet distribution for the dataset-specific component weights and the ϕ_{ij} parameter associated with the correlation between the i^{th} and j^{th} datasets.

We plot the Geweke-statistic for each chain in figure 17. No chain is perfectly behaved; as we cannot reduce to the set of stationary chains we thus exclude the most poorly behaved chains. Our lack of belief in the convergence of these

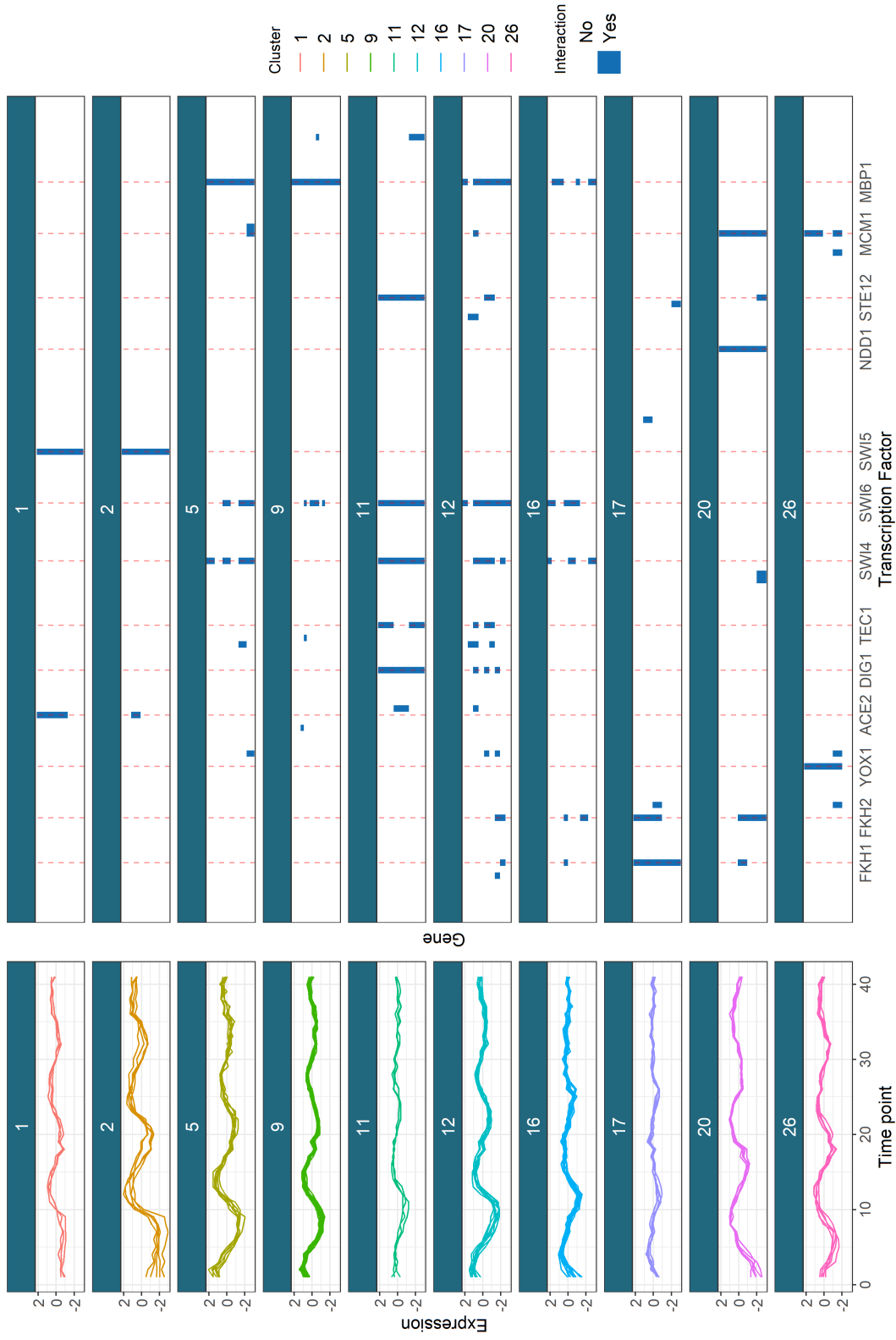


Figure 16: The clusters of genes with common labels across the time course and ChIP-chip datasets (as described in table ??). We exclude the clusters with no interactions in the ChIP-chip dataset and include a red line for the Transcription factors that dominate the clustering structure in the ChIP-chip dataset.

Within chain convergence

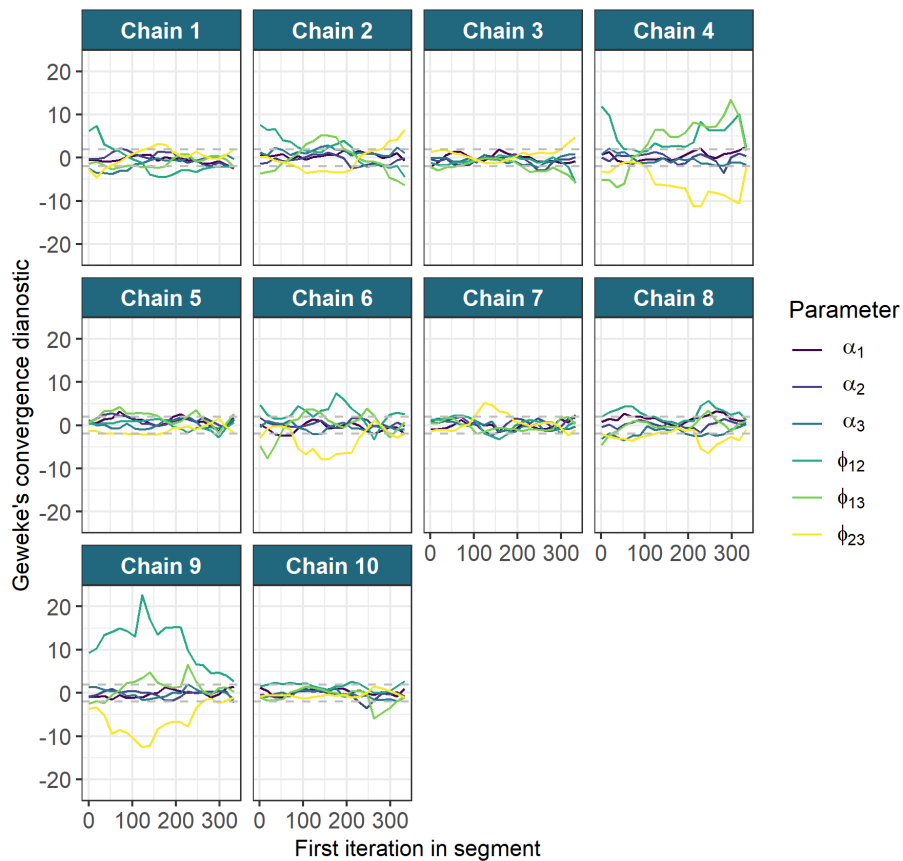


Figure 17: Chain 9 can be seen to have the most extreme behaviour in the distribution of the Geweke diagnostic for the parameters. We remove this chain from the analysis. Of the remaining chains we believe that 1, 2, 4 and 6 express the distributions furthest removed from the desired behaviour and are dropped from the analysis.

Gelman-Rubin diagnostic plot

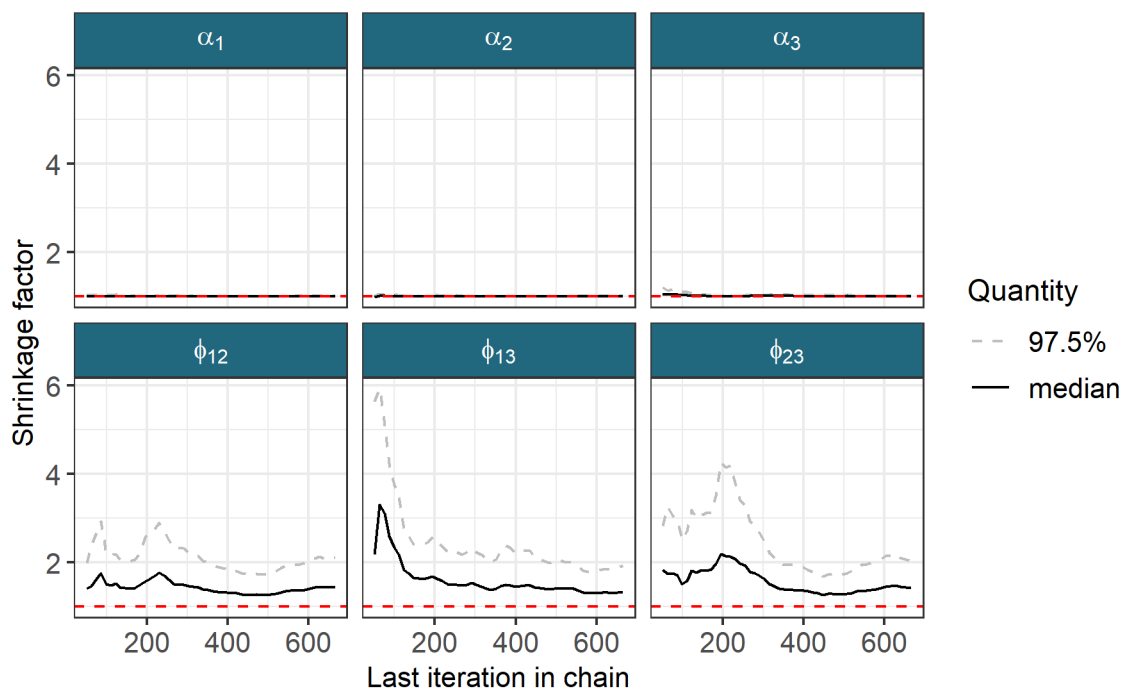


Figure 18: The chains still appear to be unconverged with \hat{R} remaining above 1.25 for the ϕ_{12} , ϕ_{13} and ϕ_{23} parameters. Stable \hat{R} is also too high with values of 1.049, 1.052 and 1.057 for ϕ_{12} , ϕ_{13} and ϕ_{23} respectively. The values of α_l cannot be seen due to the scaling of the y -axis.

chains is fortified by the behaviour of \hat{R} (which can be seen in figure 18) and the different distributions sampled for the ϕ_{lm} parameters shown in figure 19.

We visualise the the PSMs for each dataset in figure 20.

If we compare the distribution of sampled values for the ϕ parameters for the Bayesian chains that we keep based upon their convergence diagnostics, the final ensemble used ($D = 10001$, $W = 1000$) and the pooled samples from the 5 long chains, then we see that the ensemble consisting of the long chains (which might be believed to sampling different parts of the posterior distribution) is closer in its appearance to the distributions sampled by the consensus clustering than to any single chain.

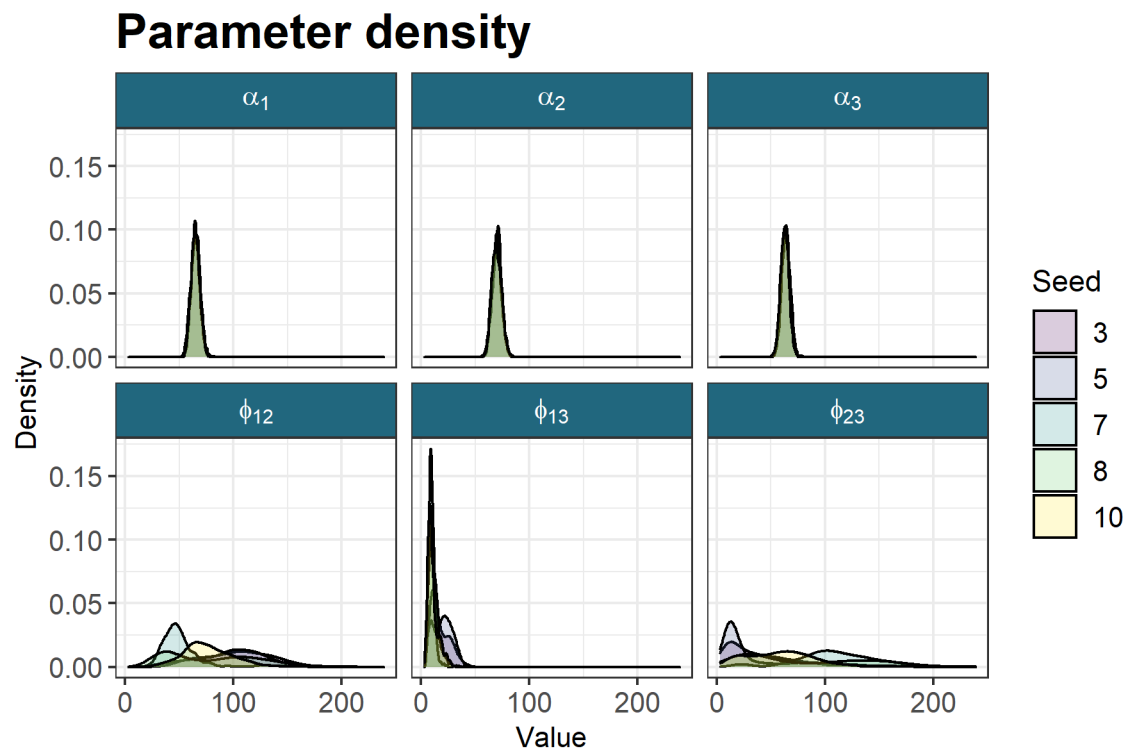


Figure 19: The densities of the continuous variables across the 5 chains kept for analysis. The mean sampled values are $\alpha_1 = 64.84$, $\alpha_2 = 69.85$, $\alpha_3 = 63.22$, $\phi_{12} = 81.76$, $\phi_{13} = 13.87$, and $\phi_{23} = 65.03$. It can be seen that different modes are being sampled for the ϕ parameters in each chain.

Multi-omics analysis

Posterior similarity matrices

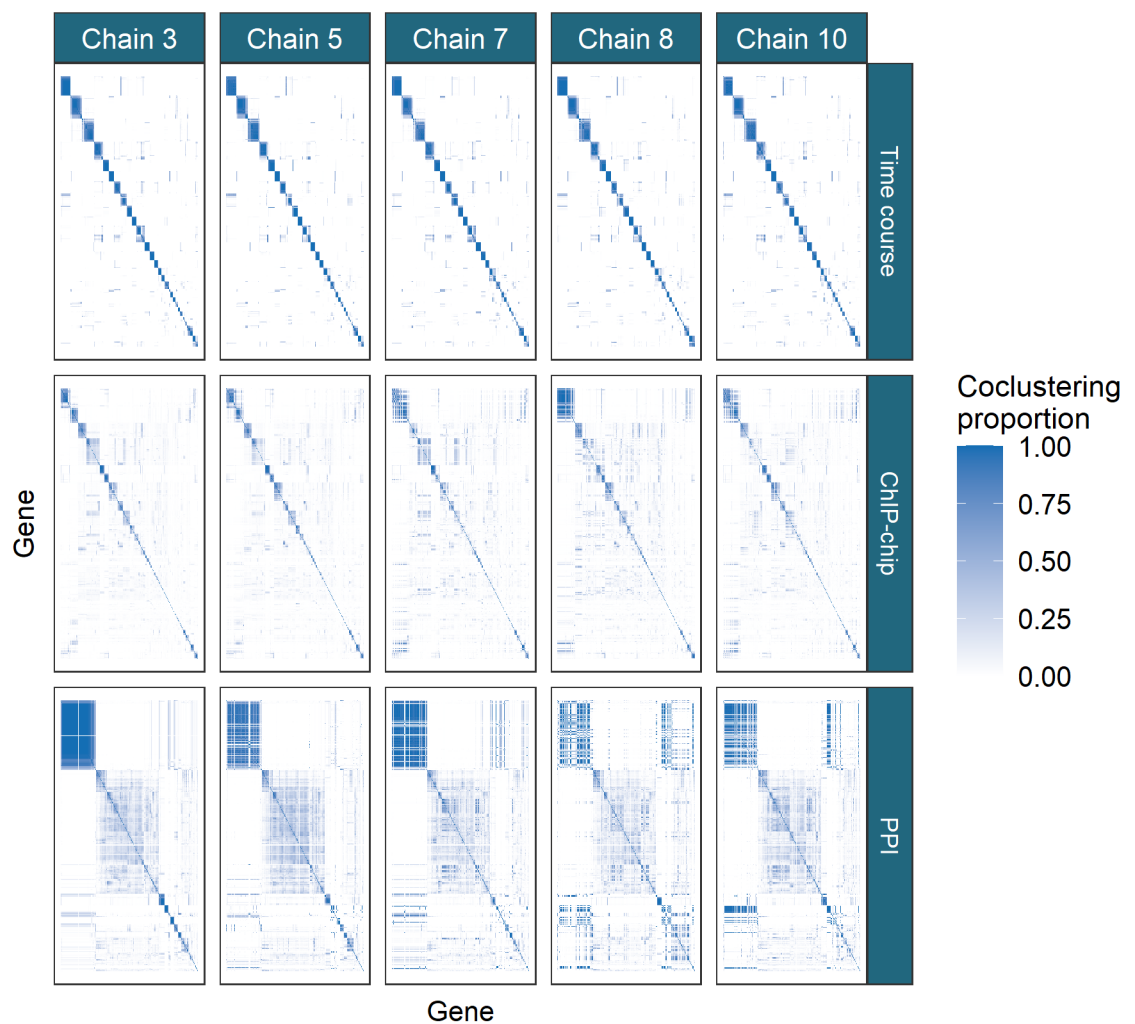


Figure 20: PSMs for each chain within each dataset. The PSMs are ordered by hierarchical clustering of the rows of the PSM for chain 3 in each dataset. There is no marked difference between the matrices for the time course data with disagreement becoming more prominent in the ChIP-chip data and more so again in the PPI dataset.

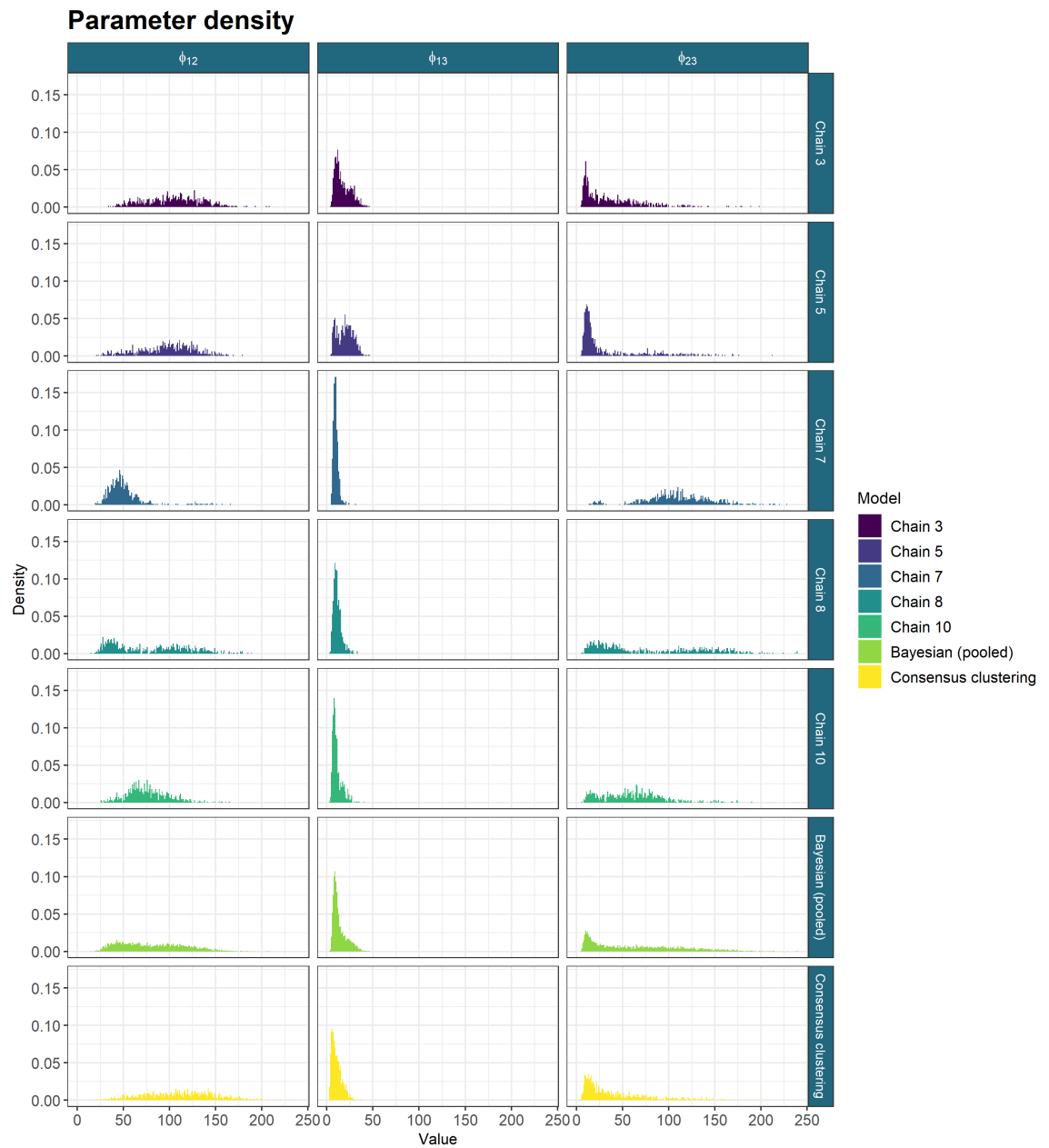


Figure 21: The sampled values for the ϕ parameters from the long chains, their pooled samples and the consensus using 1000 chains of depth 10,001. The long chains display a variety of behaviours. Across chains there is no clear consensus on the nature of the posterior distribution. The samples from any single chain are not particularly close to the behaviour of the pooled samples across all three parameters. It is the consensus clustering that most approaches this pooled behaviour.

5.3 GO term over-representation

We further show the lack of disagreement between the long chains from section 5.2 in a Gene Ontology (GO) term over-representation analysis. We estimated clusterings from the PSMs of the chains kept from section 5.2 visualised in figure 20 and the consensus matrix of the largest ensemble run (i.e. $CC(10001, 1000)$) using the `maxpear` function from the R package `mcclust` Fritsch (2012) using default settings except for `k.max` which was set to $275 \approx N/2$. To perform the GO term over-representation analysis we used the `Bioconductor` packages `clusterProfiler` (Yu et al., 2012), `biomaRt` (Durinck et al., 2009) and the annotation package `org.Sc.sgd.db` (Carlson et al., 2014).

We conditioned the test on the background set of the 551 yeast genes in the data. The gene labelled YIL167W was not found in the annotation database and was dropped from the analysis leaving a background universe of 550 genes. A hypergeometric test was used to check if the number of genes associated with specific GO terms within a cluster was greater than expected by random chance. We corrected the p -values using the Benjamini-Hochberg correction (Benjamini and Hochberg, 1995) and defined significance by a threshold of 0.01. We plotted the over-represented GO terms for the different clusterings within each dataset using the three different ontologies of “Molecular function” (**MF**), “Biological process” (**BP**) and “Cellular component” (**CC**) (figures 22, 23 and 24 respectively).

As we expect based upon the disagreement shown in figure 21, we find that the Bayesian chains have very significant disagreements between each other; there is no consensus on the results with many terms enriched in one or two chains. However, the consensus clustering finds many of the terms common to all of the long chains. This is what we would expect based upon the similarity of the ϕ_{lm} distribution in the ensemble and the pooled long chains. Consensus clustering also finds some terms with low p -values common to a majority of chains (such as DNA helicase activity in the MF ontology for the time course dataset) and a small number of GO terms unique to itself. These terms that no long chain find are normally related to other terms already over-represented within either the consensus clustering or a number of the long chains. For example, the transmembrane transporter activity and transporter activity terms uncovered by the ensemble in the time course dataset are related to terms found across 3 of the chains and by consensus clustering (specifically transferase activity and phosphotransferase).



Figure 22: GO term over-representation for the Molecular function ontology for each dataset from the final clustering of each method.



Figure 23: GO term over-representation for the Biological process ontology for each dataset from the final clustering of each method.

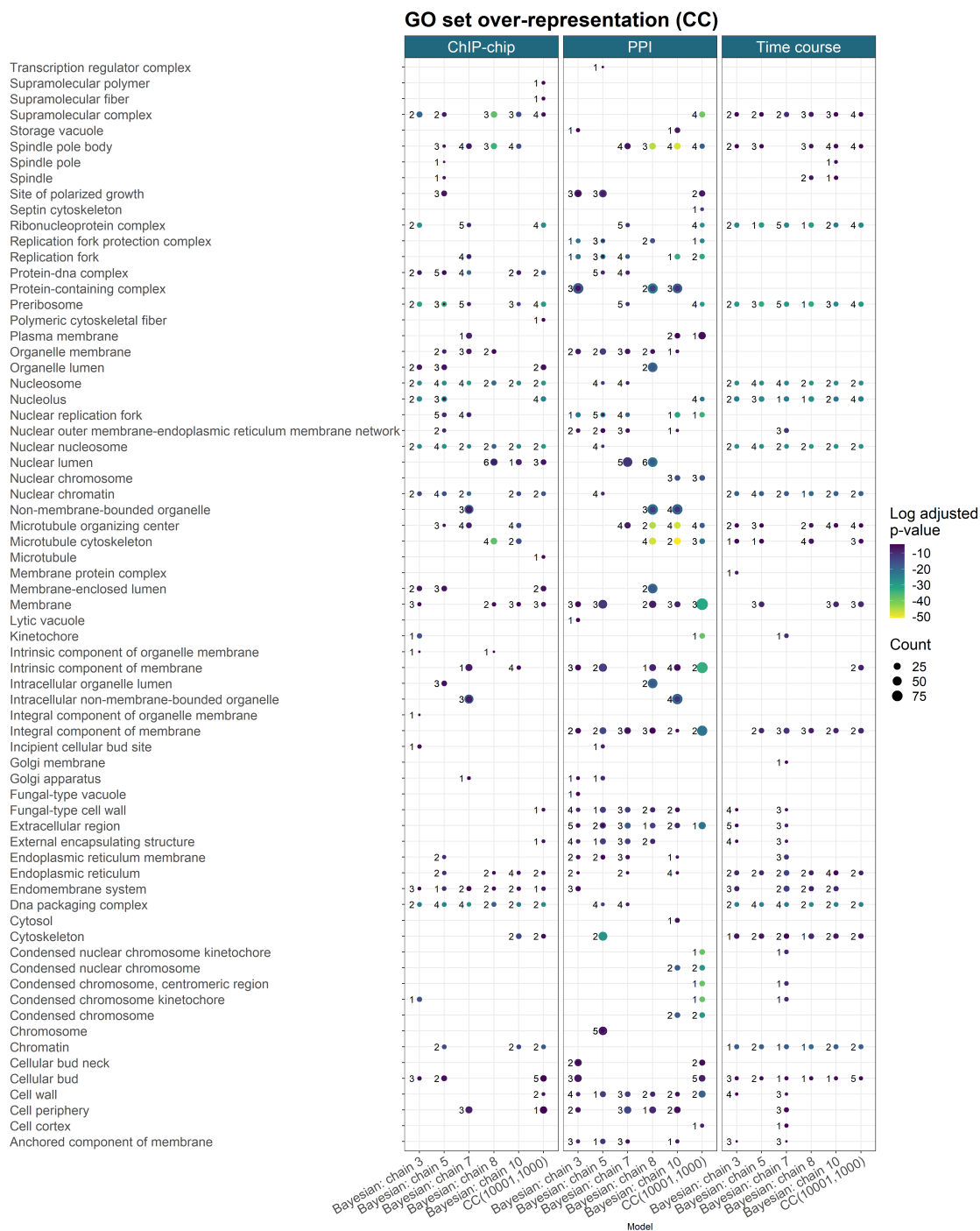


Figure 24: GO term over-representation for the Cellular component ontology for each dataset from the final clustering of each method.

References

- Sofia Aligianni, Daniel H Lackner, Steffi Klier, Gabriella Rustici, Brian T Wilhelm, Samuel Marguerat, Sandra Codlin, Alvis Brazma, Robertus AM de Bruin, and Jürg Bähler. The fission yeast homeodomain protein yox1p binds to mbf and confines mbf-dependent cell-cycle transcription to g1-s via negative feedback. *PLoS Genet*, 5(8):e1000626, 2009.
- Masashige Bando, Yuki Katou, Makiko Komata, Hirokazu Tanaka, Takehiko Itoh, Takashi Sutani, and Katsuhiko Shirahige. Csm3, tof1, and mrc1 form a heterotrimeric mediator complex that associates with dna replication forks. *Journal of Biological Chemistry*, 284(49):34355–34365, 2009.
- Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- M Carlson, S Falcon, H Pages, and N Li. Org. sc. sgd. db: Genome wide annotation for yeast. *R package version*, 2(1), 2014.
- Raymond J Cho, Michael J Campbell, Elizabeth A Winzeler, Lars Steinmetz, Andrew Conway, Lisa Wodicka, Tyra G Wolfsberg, Andrei E Gabrielian, David Landsman, David J Lockhart, et al. A genome-wide transcriptional analysis of the mitotic cell cycle. *Molecular cell*, 2(1):65–73, 1998.
- Steffen Durinck, Paul T Spellman, Ewan Birney, and Wolfgang Huber. Mapping identifiers for the integration of genomic datasets with the r/bioconductor package biomaRt. *Nature protocols*, 4(8):1184, 2009.
- Yasin Şenbabaoğlu, George Michailidis, and Jun Z Li. A reassessment of consensus clustering for class discovery. *bioRxiv*, page 002642, 2014a.
- Yasin Şenbabaoğlu, George Michailidis, and Jun Z Li. Critical limitations of consensus clustering in class discovery. *Scientific reports*, 4(1):1–13, 2014b.
- Mark E Ewen. Where the cell cycle and histones meet. *Genes & development*, 14(18):2265–2270, 2000.
- Thomas S Ferguson. A bayesian analysis of some nonparametric problems. *The annals of statistics*, pages 209–230, 1973.
- Arno Fritsch. *mcclust: Process an MCMC Sample of Clusterings*, 2012. URL <https://CRAN.R-project.org/package=mcclust>. R package version 1.0.
- Arno Fritsch, Katja Ickstadt, et al. Improved criteria for clustering based on the posterior similarity matrix. *Bayesian analysis*, 4(2):367–391, 2009.
- Evelina Gabasova, John Reid, and Lorenz Wernisch. Clusternomics: Integrative context-dependent clustering for heterogeneous datasets. *PLoS computational biology*, 13(10):e1005781, 2017.

- Andrew Gelman, Donald B Rubin, et al. Inference from iterative simulation using multiple sequences. *Statistical science*, 7(4):457–472, 1992.
- Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC press, 2013.
- John Geweke et al. *Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments*, volume 196. Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN, 1991.
- Vishwanath R Iyer, Christine E Horak, Charles S Scafe, David Botstein, Michael Snyder, and Patrick O Brown. Genomic binding sites of the yeast cell-cycle transcription factors sbf and mbf. *Nature*, 409(6819):533–538, 2001.
- Paul Kirk, Jim E Griffin, Richard S Savage, Zoubin Ghahramani, and David L Wild. Bayesian correlated clustering to integrate multiple datasets. *Bioinformatics*, 28(24):3290–3297, 2012.
- Christina Knudson and Dootika Vats. *stableGR: A Stable Gelman-Rubin Diagnostic for Markov Chain Monte Carlo*, 2020. URL <https://CRAN.R-project.org/package=stableGR>. R package version 1.0.
- Manfred Koranda, Alexander Schleiffer, Lukas Endler, and Gustav Ammerer. Forkhead-like transcription factors recruit ndd1 to the chromatin of g2/m-specific promoters. *Nature*, 406(6791):94–98, 2000.
- Raman Kumar, David M Reynolds, Andrej Shevchenko, Anna Shevchenko, Sherilyn D Goldstone, and Stephen Dalton. Forkhead transcription factors, fkh1p and fkh2p, collaborate with mcm1p to control transcription required for m-phase. *Current Biology*, 10(15):896–906, 2000.
- Jessica P Lao, Katie M Ulrich, Jeffrey R Johnson, Billy W Newton, Ajay A Vashisht, James A Wohlschlegel, Nevan J Krogan, and David P Toczyski. The yeast dna damage checkpoint kinase rad53 targets the exoribonuclease, xrn1. *G3: Genes, Genomes, Genetics*, 8(12):3931–3944, 2018.
- Samuel A Mason, Faiz Sayyid, Paul DW Kirk, Colin Starr, and David L Wild. Mdi-gpu: accelerating integrative modelling for genomic-scale data using gpu computing. *Statistical Applications in Genetics and Molecular Biology*, 15(1):83–86, 2016.
- Helen J McBride, Yaxin Yu, and David J Stillman. Distinct regions of the swi5 and ace2 transcription factors are required for specific gene activation. *Journal of Biological Chemistry*, 274(30):21029–21036, 1999.
- Stefano Monti, Pablo Tamayo, Jill Mesirov, and Todd Golub. Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. *Machine learning*, 52(1-2):91–118, 2003.

- Tata Pramila, Shawna Miles, Debraj GuhaThakurta, Dave Jemiolo, and Linda L Breeden. Conserved homeodomain proteins interact with mads box protein mcm1 to restrict ecb-dependent transcription to the m/g1 phase of the cell cycle. *Genes & development*, 16(23):3034–3045, 2002.
- Sylvia Richardson and Peter J Green. On bayesian analysis of mixtures with an unknown number of components. *Journal of the Royal Statistical Society: series B*, 59(4):731–792, 1997.
- Judith Rousseau and Kerrie Mengersen. Asymptotic behaviour of the posterior distribution in overfitted mixture models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(5):689–710, 2011.
- Richard S Savage, Zoubin Ghahramani, Jim E Griffin, Bernard J De la Cruz, and David L Wild. Discovering transcriptional modules by bayesian data integration. *Bioinformatics*, 26(12):i158–i167, 2010.
- Gideon Schwarz et al. Estimating the dimension of a model. *The annals of statistics*, 6(2):461–464, 1978.
- Samuel Sanford Shapiro and Martin B Wilk. An analysis of variance test for normality (complete samples). *Biometrika*, 52(3/4):591–611, 1965.
- Itamar Simon, John Barnett, Nancy Hannett, Christopher T Harbison, Nicola J Rinaldi, Thomas L Volkert, John J Wyrick, Julia Zeitlinger, David K Gifford, Tommi S Jaakkola, et al. Serial regulation of transcriptional regulators in the yeast cell cycle. *Cell*, 106(6):697–708, 2001.
- Paul T Spellman, Gavin Sherlock, Michael Q Zhang, Vishwanath R Iyer, Kirk Anders, Michael B Eisen, Patrick O Brown, David Botstein, and Bruce Futcher. Comprehensive identification of cell cycle-regulated genes of the yeast *saccharomyces cerevisiae* by microarray hybridization. *Molecular biology of the cell*, 9(12):3273–3297, 1998.
- Julia van der Felden, Sarah Weisser, Stefan Brückner, Peter Lenz, and Hans-Ulrich Mösch. The transcription factors *tec1* and *ste12* interact with coregulators *msa1* and *msa2* to activate adhesion and multicellular development. *Molecular and cellular biology*, 34(12):2283–2293, 2014.
- Zoé Van Havre, Nicole White, Judith Rousseau, and Kerrie Mengersen. Overfitting bayesian mixture models with an unknown number of components. *PloS one*, 10(7):e0131739, 2015.
- Dootika Vats and Christina Knudson. Revisiting the gelman-rubin diagnostic. *arXiv preprint arXiv:1812.09384*, 2018.
- Guangchuang Yu, Li-Gen Wang, Yanyan Han, and Qing-Yu He. clusterProfiler: an r package for comparing biological themes among gene clusters. *OMICS: A Journal of Integrative Biology*, 16(5):284–287, 2012. doi: 10.1089/omi.2011.0118.