

Bayesian model-based clustering for populations of network data

Anastasia Mantziou¹, Simón Lunagómez², and Robin Mitra³

¹The Alan Turing Institute

²Departamento de Estadística, ITAM

³Department of Statistical Science, University College London

Abstract

There is increasing appetite for analysing populations of network data due to the fast-growing body of applications demanding such methods. While methods exist to provide readily interpretable summaries of heterogeneous network populations, these are often descriptive or ad hoc, lacking any formal justification. In contrast, principled analysis methods often provide results difficult to relate back to the applied problem of interest. Motivated by two complementary applied examples, we develop a Bayesian framework to appropriately model complex heterogeneous network populations, whilst also allowing analysts to gain insights from the data, and make inferences most relevant to their needs. The first application involves a study in Computer Science measuring human movements across a University. The second analyses data from Neuroscience investigating relationships between different regions of the brain. While both applications entail analysis of a heterogeneous population of networks, network sizes vary considerably. We focus on the problem of clustering the elements of a network population, where each cluster is characterised by a network representative. We take advantage of the Bayesian machinery to simultaneously infer the cluster membership, the representatives, and the community structure of the representatives, thus allowing intuitive inferences to be made. The implementation of our method on the human movement study reveals interesting movement patterns of individuals in clusters, readily characterised by their network representative. For the brain networks application, our model reveals a cluster of individuals with different network properties of particular interest in Neuroscience. The performance of our method is additionally validated in extensive simulation studies.

Keywords: Bayesian models, Clustering, Mixture models, Populations of network data, Object data analysis

1 Introduction

Conventional statistical methods for modelling and analysing data, such as standard regression models, assume each observation or outcome value is a scalar. While this is appropriate for many applications there is sometimes a need to handle more complex types of data. For example, each observation may constitute a set of interconnected points where the distribution of the connections between the points determines the observation's properties. Such data are typically referred to as networks in the literature, and modelling them is fundamentally important in some applications. For example, suppose we are interested in monitoring movements of subjects across a geographical area via displays at fixed locations. Each display represents a point, or node, of the network and when a subject moves from one display to another a connection, or edge, between the two nodes is assumed. The pattern of movement across the displays then characterises a network for that subject. Effectively modelling such data can provide analysts with important insights into problems of interest. For example, we may be interested in distinguishing different patterns of behaviour among the different subjects. Are some subjects more likely to visit certain displays than others? Do subjects have different or unusual patterns of movements to the majority of subjects? These are all important applied questions that present significant challenges to address due to the complexity of dealing with network data.

The availability of populations of network data has risen substantially in recent years, due to the advancement of technological means that record this type of data (White et al., 1986; Fields and Song, 1989). This has inspired many researchers to develop statistical models that most accurately describe the probabilistic mechanism that generates a network population. Specifically, there are three different frameworks considered in the literature for modelling populations of network data: 1) the latent space models, 2) the distance-based models, and 3) the measurement error models.

For a single network observation, the fundamental idea behind the latent space class of models is that the occurrence of an edge between two nodes depends on the positions of the nodes in a latent space (Hoff et al., 2002; Young and Scheinerman, 2007). Recent studies (Gollini and Murphy, 2016; Levin et al., 2017; Durante et al., 2017; Wang et al., 2019; Nielsen and Witten, 2018; Arroyo et al., 2021) on modelling populations of network data have extended this idea to build models for populations of networks with aligned vertex sets, assuming that the nodes lie in a common, unobserved subspace.

Another approach to modelling populations of network data is the utilisation of distance metrics that measure similarities among networks with respect to global or local characteristics of the networks (Donnat and Holmes, 2018). Under this framework, researchers rely on the notion of an average network that represents a network population, with respect to a specified distance metric (Lunagómez et al., 2021; Kolaczyk et al., 2020; Ginestet et al., 2017).

The third class of models, measurement error models, account for the erroneous nature of the networks. A fundamental source of noise found in network data originates from the various measurement tools used for the construction of networks, i.e. the processes used to measure an interaction (edge) between two objects (nodes). Researchers focusing on the statistical analysis of networks as single observations have developed methods to incorporate the uncertainty of falsely observing edges or non-edges in a network. Such studies involve predicting network topologies accounting for the falsely non-observed edges (Jiang et al., 2011), estimating the adjacency matrix from a set of noisy entries (Chatterjee, 2015), classifying nodes of networks with errorful edges (Priebe et al., 2015), developing a regression model for networks assuming that the observed network is a perturbed version of a true unobserved network (Le and Li, 2020) and performing Bayesian inference on the network’s structure utilising information from measurements (Young et al., 2020). Another group of studies focuses on the propagation of the error to network summary statistics (Balachandran et al., 2017; Chang et al., 2020), and to estimators of average causal effects under network interference when the error arises from a measurement process used to construct the network (Li et al., 2021). Le et al. (2018), Newman (2018) and Peixoto (2018) develop this idea to model populations of network observations in order to infer the probabilistic mechanism that generates the network population. Specifically they assume the networks are noisy realisations of a true unobserved network. Similarly, Josephs et al. (2021) consider the problem of network recovery from multiple noisy network realisations, for unlabeled networks.

Despite the growing research interest on modelling populations of network data, only few studies developed to date consider the heterogeneity that can exist in a network population. Notably, Mukherjee et al. (2017) were the first to consider the problem of clustering populations of network data. They assume two different scenarios, (a) the networks in the population share the same set of nodes, and (b) the networks in the population do not share the same set of nodes. Our paper focuses on Scenario (a). In Scenario (a), the authors obtain a mixture model of graphons and implement a spectral clustering algorithm to infer the membership allocation of each network observation.

An application driven study on clustering populations of network data is introduced by Diquigiovanni and Scarpa (2019), who aim to cluster a population of networks where each network observation represents the playing style of a football team at a specific match. The clustering approach seen in this study involves the specification of an ad hoc measure of similarity between networks, and the implementation of an agglomerative method for clustering the networks according to their similarities.

To the best of our knowledge, the third and last study that examines the problem of clustering network populations is that of Signorelli and Wit (2020). In this study, the authors deal with the problem of clustering using a mixture model whose components can be any statistical network model, under the restriction that it can be specified as a Generalised Linear Model (GLM). For estimating the parameters of their model, they implement an Expectation Maximisation (EM) algorithm for a predefined number of clusters. To determine the network model for their mixture, the authors propose the initial use of a mixture of saturated network models, to reveal information about the structure of the data at hand. The saturated network model assumes that each edge in each network in the population is generated with some unique, unconstrained probability.

Another group of studies that accounts for the heterogeneity in a set of network observations are the studies that perform the task of network classification. Some of these studies consider either specific network summary measures (Prasad et al., 2015), or vectorise only the important entries (edges) of the adjacency matrix (Richiardi et al., 2011; Zhang et al., 2012) to classify networks. Thus, they ignore the overall networks’ structure. In contrast to these studies, Arroyo Relión et al. (2019) perform prediction of the class membership of networks using a linear classifier with the adjacency matrices of the networks as predictors. Their approach accounts for the networks’ structure by using a penalty to select important nodes and edges.

These contributions provide interesting approaches for identifying variations between network data, but

there are some key limitations associated with these:

- In the studies of [Mukherjee et al. \(2017\)](#), [Diquigiovanni and Scarpa \(2019\)](#) and [Arroyo Reli3n et al. \(2019\)](#), the methods proposed are non model-based. [Mukherjee et al. \(2017\)](#) and [Diquigiovanni and Scarpa \(2019\)](#) propose algorithms that detect underlying network clusters in the data and [Arroyo Reli3n et al. \(2019\)](#) predict class membership of the networks. In all three studies the groups of networks identified cannot be interpreted using a parametric representation. Interpretability of the different groups of networks in a population is crucial in many applications in order to infer group specific properties and differences.
- While [Signorelli and Wit \(2020\)](#) provide a model-based approach for clustering populations of network data, the mixture components must conform to rigid modelling assumptions. This means that only specific characteristics of the networks can be inferred depending on what these model assumptions allow. It would be ideal to have a framework that is flexible enough to incorporate different modelling assumptions as deemed appropriate to application allowing the most scientifically relevant inferences to be made.
- In addition, [Signorelli and Wit \(2020\)](#) propose to initially obtain a mixture of saturated network models, thus resulting in an overly complex model with a large number of parameters to estimate. This can substantially increase the computational time needed for the EM algorithm to converge as well as increasing the potential for non-convergence due to having to explore a very high dimensional parameter space.
- The supervised approach of [Arroyo Reli3n et al. \(2019\)](#) requires a training data set to predict the class. The class labels of the networks in the training data set must be pre-specified, which can be restrictive for some network applications for which we do not have a priori information about the networks.

To address these limitations, in this paper we propose a mixture model for identifying clusters of networks in a network population, with respect to similarities detected in the connectivity patterns of the networks' nodes. We consider the case when the networks in the population share the same set of n nodes, and each network could belong to one of, a predefined number of, C clusters. Inspired by the approach of [Le et al. \(2018\)](#), we adopt a measurement error formulation, assuming networks lying within each of the C clusters are noisy realisations of a true underlying network representative. The attractive feature of this approach is that it decouples the statistical model for the network data from the underlying cluster specific network properties. We are thus able to provide a flexible model-based approach for detecting clusters of networks in a network population, as well as interpret these clusters with respect to our model parameterisation. Our framework is also flexible enough to incorporate, and thus exploit, any underlying assumptions about the structure of the networks within the clusters that are of scientific interest or otherwise supported by the data.

[Le et al. \(2018\)](#) develop a model for populations of network observations assuming that noisy network-valued observations arise from a true underlying adjacency matrix. The inferential framework built in their study consists of two steps. First, they use a Spectral Clustering algorithm to infer the community structure formed by the nodes of the true underlying network, and second they implement an EM algorithm to estimate the model parameters. An evident limitation of their inferential framework is that their algorithm does not simultaneously update the parameters of their model for the network data and the parameters characterising the underlying network structure, as this would require the development of new techniques. In addition, the assumption of a sole true underlying adjacency matrix is quite restrictive, especially for a large sample of networks where a degree of heterogeneity is expected.

We adopt a Bayesian modelling approach which provides some unique advantages over previous approaches. In particular by utilising Markov Chain Monte Carlo (MCMC) methods, we are able to simultaneously infer the cluster membership of the networks, together with model parameters characterising the distribution of the networks within each cluster as well as those that characterise the structure of the underlying cluster specific network representatives. To best of our knowledge, there is no coherent framework in the literature that permits this type of complete inference from the network data. Our framework is flexible enough to answer a diverse range of applied questions with respect to the heterogeneity in a network population. These include being able to detect clusters of networks as well as inferring key different features between clusters through comparisons between the underlying representatives. In addition, interest may lie in identifying observations that do not follow the distribution of the majority of the network data, and the framework can also be formulated to detect outlying network observations.

Our approach is motivated by two applied examples in very different fields, one involving monitoring movement of people across a University Campus and another measuring individuals’ connectivity patterns across different regions of the brain. In this second example we are particularly interested in determining whether any individuals possess unusual connectivity patterns, and if so how these differ to the rest of the sample. This applied question can be readily addressed using the proposed framework with the outlier formulation described above. Other principled methods, by contrast, would struggle to produce readily interpretable summaries.

The remainder of this article is organised as follows. In Section 2, we describe the applied examples that motivated the development of the methods. In Section 3, we provide background to our modelling framework. In Section 4, we develop the Bayesian formulation of a mixture of measurement error models for network data, along with the MCMC scheme to make inferences. We further present the Sparse Finite Mixture (SFM) extension of our model that allows inferences for the number of clusters. In Section 5, we present simulation studies to assess the performance of our method for various network sizes and sample sizes. In Section 6, we analyse the two different motivating populations of network data examples to provide important insights into applied questions of interest. This also serves to demonstrate the broad applicability of our methods. Lastly, in Section 7, we give some concluding remarks.

2 Two motivating examples

In this section we introduce two different applied examples that have triggered research questions which we aim to answer in our study. While both applications address very different areas, a common feature in both data sets is the heterogeneity of the network data in the sample.

2.1 Data on movements of subjects across a University Campus

The first example comes from data collected on movement of people across Lancaster University Campus in the UK. The study was performed by members of the Computing and Communications department at the university (Shaw et al., 2018).

A series of fixed displays are located across the campus (Figure 1 left). Individuals taking part in the study installed a Tacita mobile application on their phone, and whenever they pass one of these displays this application registers their presence at that location. The application can also serve as a means of communication between a display and a viewer, in order for the viewers to be able to see content relevant to their interests. Specifically, the viewer can request what to see on the screen of the display, but also the display can detect when a user is in its proximity in order to show content aligned with their interests. Thus, the application records the consecutive displays visited by the users, along with the time visited, and the type of content shown.

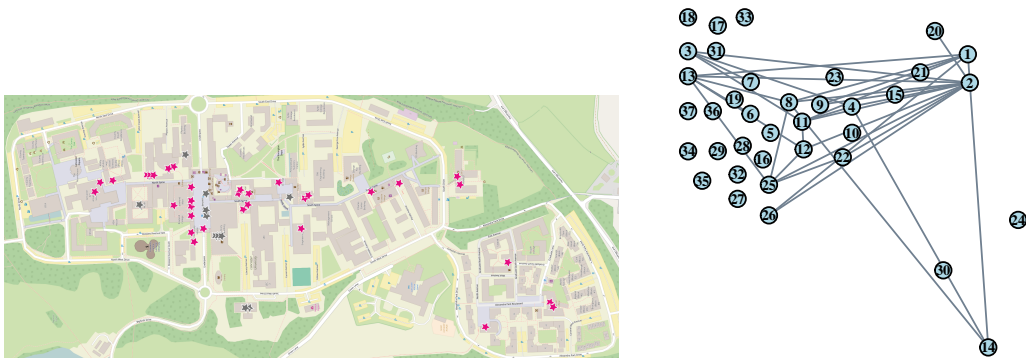


Figure 1: Left: Lancaster University campus map with stars indicating displays’ location. Right: Network visualisation of movements of one individual from the sample. The network layout corresponds to the physical location of the displays.

Consequently, the Tacita data set serves as an example of a population of network observations, as the movements of each individual can be represented by a network. Thus, for each individual we can obtain a network where nodes represent displays, and edges represent the movements of the user among the displays.

In this example each network corresponds to the aggregated movements of an individual across the campus during different times of the day, resulting in 120 network observations sharing the same set of 37 nodes. In Figure 1 (right), we illustrate the movements among the 37 displays (nodes) on campus, for one of the users of the Tacita application. The nodes' positions in Figure 1 correspond to the physical location of the displays. We note here that there is not direct correspondence between the display locations indicated on the campus map and the nodes positions in Figure 1 for two reasons: first, after data manipulation and consultation with our collaborators, some displays were not considered in our analysis, and second, the data collected involved newly activated displays not depicted on the campus map in Figure 1. The following questions arise:

- Can we infer a meaningful number of clusters based on the observed population of individuals?
- Can we detect different patterns among the users' movements?
- Can we cluster the users according to their movements?
- How informative can the clustering be for the users in our data?

2.2 Data measuring connectivity patterns across different regions of the brain

The second example is a population of networks data set arising from the field of Neuroscience. In this data example, connectivity patterns across different regions of the brain were measured for 30 healthy individuals at resting state. For each individual a series of 10 measurements were taken using diffusion magnetic resonance imaging (dMRI). The measurements are represented as networks with the nodes corresponding to fixed regions of the brain, and edges denoting the connections recorded among those regions. Specifically, the network data consist of 200 nodes (regions of the brain) according to the CC200 atlas (Craddock et al., 2012), and the resulting data set consists of 300 network observations. In Figure 2, we illustrate the network representation of one brain scan for one of the individuals in the data set.

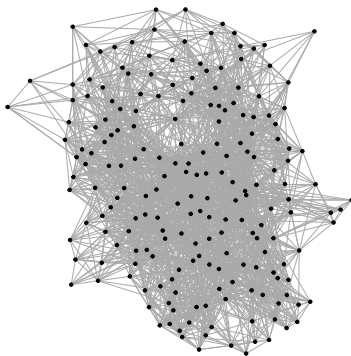


Figure 2: Network representation of brain scan taken for one of the individuals in the sample.

This data set has been also discussed in the study of Zuo et al. (2014), Arroyo et al. (2021) and Lunagómez et al. (2021), with the latter two studies analysing the data from a modelling perspective. Specifically, Arroyo et al. (2021) investigate the ability of their method to identify differences among individuals with respect to communities formed from the networks' nodes, while Lunagómez et al. (2021) assume unimodality of the probabilistic mechanism that generates the network population and infer a representative network for the population of individuals, according to a pre-specified distance metric. None of these approaches seek to determine and interpret clusters of networks. In particular, a relevant research question here is determining whether any individuals meaningfully differ from the rest of the sample in terms of their network characteristics, and if so, in what way. Motivated by this objective, we can formulate our model to capture, and subsequently interpret, outlying network data, through an appropriate cluster specification, as presented in Section 4.4.

More generally, our goal in this application is to explore possible heterogeneity amongst the networks. Questions include:

- Can we identify clusters of individuals with respect to similarities found in their connectivity patterns? In particular, can we identify individuals with different network characteristics to the majority of the population?

- Can we interpret the clusters identified with respect to some network feature so they are relevant to Neuroscience applications?
- Are brain scans of the same individual assigned to the same cluster?

These applied research questions have motivated the proposed mixture of network measurement error models, described in detail in Section 6.

3 Background

A network can be represented as a graph $\mathcal{G} = (V, E)$, where $V = \{1, \dots, n\}$ represents the set of n nodes and E represents the set of observed edges in \mathcal{G} , with $E \in \mathcal{E}_n$ and $\mathcal{E}_n = \{(i, j) | i, j \in V\}$. A common mathematical network representation is an $n \times n$ adjacency matrix $A_{\mathcal{G}}$, with the $A_{\mathcal{G}}(i, j)$ entry of the matrix denoting the state of the (i, j) edge. The $(i, j)^{th}$ element of the adjacency matrix for a graph with binary edges is,

$$A_{\mathcal{G}}(i, j) = \begin{cases} 1, & \text{if an edge occurs between nodes } i \text{ and } j, \\ 0, & \text{otherwise.} \end{cases}$$

The adjacency matrix of an undirected graph with no self-loops is symmetric with $A_{\mathcal{G}}(i, j) = A_{\mathcal{G}}(j, i)$ and $A_{\mathcal{G}}(i, i) = 0$, for $i, j \in \{1, \dots, n\}$. By $\mathcal{G}_1, \dots, \mathcal{G}_N$ we represent a population of N graphs, with corresponding adjacency matrices $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$. In this study, we assume that the networks in the population $\mathcal{G}_1, \dots, \mathcal{G}_N$ are undirected with no self-loops, and share the same set of n nodes. We note here that our methods can be equally applied to populations of directed graphs.

When modelling a population of network observations it is often natural to assume that the networks have been subjected to some noise or measurement error during their construction. For example, in the application investigating the movement of subjects across campus, it is possible that the Tacita mobile application fails to register a subject at a display occasionally, while also sometimes incorrectly registering a subject at a display, particularly when displays are located fairly close together. This results in network data that might have edges missing as well as edges recorded that should not be present.

Under this measurement error hypothesis, the researcher assumes that the observed network data correspond to noisy realisations of a true underlying network, which leads to recording some erroneous edges, due to the existence of an underlying measurement error process. [Le et al. \(2018\)](#) were the first to introduce this approach for modelling populations of networks, which has inspired the proposed modelling approach.

[Le et al. \(2018\)](#) assume that the information contained in the network population can be summarised by a representative network \mathcal{G}^* , and a measurement error process that does not allow us to accurately observe the representative network. Specifically, the authors assume a false positive probability P_{ij} of observing an edge between nodes i, j in the k^{th} network observation \mathcal{G}_k , given that there is no edge between the same two nodes in the representative network \mathcal{G}^* ; and respectively, a false negative probability Q_{ij} of not observing an edge for the nodes i, j in the k^{th} network observation \mathcal{G}_k , while there is an edge for the same two nodes in the representative network \mathcal{G}^* . Thus, the entries of the matrices P, Q , are the false positive/negative probabilities of seeing/not seeing an edge respectively between two nodes in the data.

The mathematical formulation of the above set-up is the following. Let $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$ denote the adjacency matrices for the network population, and $A_{\mathcal{G}^*}$ the adjacency matrix of the representative network. The false positive and false negative probabilities P_{ij}, Q_{ij} can be described as follows,

$$\begin{aligned} \text{if } A_{\mathcal{G}^*}(i, j) = 1, \text{ then } A_{\mathcal{G}_k}(i, j) &= \begin{cases} 1, & \text{with prob } 1 - Q_{ij} \\ 0, & \text{with prob } Q_{ij} \end{cases} \quad ; \\ \text{if } A_{\mathcal{G}^*}(i, j) = 0, \text{ then } A_{\mathcal{G}_k}(i, j) &= \begin{cases} 1, & \text{with prob } P_{ij} \\ 0, & \text{with prob } 1 - P_{ij} \end{cases} \quad . \end{aligned} \quad (1)$$

From (1) it follows that the probability of the occurrence or non-occurrence of an edge between nodes i, j in the k^{th} network observation is,

$$P(A_{\mathcal{G}_k}(i, j) | A_{\mathcal{G}^*}(i, j) = 1, P_{ij}, Q_{ij}) = (1 - Q_{ij})^{A_{\mathcal{G}_k}(i, j)} \cdot Q_{ij}^{1 - A_{\mathcal{G}_k}(i, j)}, \text{ if } A_{\mathcal{G}^*}(i, j) = 1;$$

$$P(A_{\mathcal{G}_k}(i, j) | A_{\mathcal{G}^*}(i, j) = 0, P_{ij}, Q_{ij}) = P_{ij}^{A_{\mathcal{G}_k}(i, j)} \cdot (1 - P_{ij})^{1 - A_{\mathcal{G}_k}(i, j)}, \text{ if } A_{\mathcal{G}^*}(i, j) = 0.$$

Le et al. (2018) treat the adjacency matrix of the representative network $A_{\mathcal{G}^*}$ as a latent variable, while the false positive and false negative probabilities P_{ij} and Q_{ij} are model parameters. Thus, the likelihood of the representative network $A_{\mathcal{G}^*}$ given the network data $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$, as seen in Le et al. (2018), is

$$\begin{aligned} \mathcal{L}(A_{\mathcal{G}^*}; A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) &= \prod_{k=1}^N \prod_{(i, j): i < j} [(1 - Q_{ij})^{A_{\mathcal{G}_k}(i, j)} \cdot Q_{ij}^{1 - A_{\mathcal{G}_k}(i, j)} \cdot W_{ij}]^{A_{\mathcal{G}^*}(i, j)} \\ &\quad [P_{ij}^{A_{\mathcal{G}_k}(i, j)} \cdot (1 - P_{ij})^{1 - A_{\mathcal{G}_k}(i, j)} \cdot (1 - W_{ij})]^{1 - A_{\mathcal{G}^*}(i, j)}, \end{aligned}$$

where $W_{ij} = \mathbb{E}A_{\mathcal{G}^*}(i, j)$ represents the probability of observing an edge between nodes i, j in $A_{\mathcal{G}^*}$.

Le et al. (2018) further assume that the nodes of the underlying true network form communities that can be described by a Stochastic Block Model (SBM). Under the SBM assumption, each node of the true network belongs to an unobserved block $k \in \{1, \dots, K\}$, and the probability of observing an edge between two nodes depends on their block membership denoted by $\{b_i\}_{i=1}^n$, with $b_i \in \{1, \dots, K\}$. The probability of observing an edge between nodes (i, j) with $b_i = k, b_j = l$ is represented by θ_{kl} . In addition, the corresponding block structure is assumed to be shared among the matrices P, Q and W .

The inference of the model parameters and the latent variable is conducted in two stages. First, a Spectral Clustering algorithm is applied to reveal the underlying block membership of the representative's nodes, and second, an EM algorithm is implemented to estimate the model parameters. While this formulation has appealing features it would be ideal to have a coherent modelling framework that can jointly infer block membership of the representative's nodes together with the parameters characterising the distribution of the network data. In addition, using an EM algorithm to estimate model parameters means that measures of uncertainty such as standard errors rely on asymptotic approximations that may not be valid in many applications, particularly when involving small samples sizes.

In the next section we propose a mixture of measurement error models inspired by Le et al. (2018) for clustering heterogeneous network data. We adopt a Bayesian framework that allows us to jointly infer the parameters of the measurement error model as well as those characterising the underlying network representatives corresponding to each cluster. In addition, the Bayesian formulation is flexible enough to accommodate diverse modelling assumptions for the network representatives.

4 A Mixture of Measurement Error models

In this section, we detail the formulation and implementation of the mixture of measurement error models. We first describe the Bayesian formulation of the measurement error model when there is only one cluster. We then extend this to multiple clusters. Following this we describe how posterior samples can be obtained using MCMC. Finally we describe a special case of this formulation that can correspond to detecting outlying networks in the data.

4.1 Model formulation

To begin with, we assume underlying the network data there is a latent representative network with adjacency matrix denoted by $A_{\mathcal{G}^*}$. In addition, we assume that the probability of observing a false positive or false negative edge between two nodes in the network data is independent of the pair of nodes considered. Thus, the false positive probability, p , and false negative probability, q , can be viewed as scalars. In this way, we limit model complexity in terms of the number of parameters to be inferred. The specification of component specific $n \times n$ matrices of false positive probabilities and false negative probabilities would lead to a drastic increase in the number of model parameters. A compromise assumes matrices P, Q share an SBM structure defined by the true network $A_{\mathcal{G}^*}$, as in Le et al. (2018). However this assumption would only be relevant where an SBM structure is already known to be appropriate, which might not be appropriate for some applications. We discuss this further in Section 6.

Under this specification, the probability mass function of the edge state between nodes (i, j) in network observation k , given $A_{\mathcal{G}^*}(i, j), p, q$, is

$$\begin{aligned} P(A_{\mathcal{G}_k}(i, j) | A_{\mathcal{G}^*}(i, j), p, q) &= [(1 - q)^{A_{\mathcal{G}_k}(i, j)} \cdot q^{1 - A_{\mathcal{G}_k}(i, j)}]^{A_{\mathcal{G}^*}(i, j)} \\ &\quad [p^{A_{\mathcal{G}_k}(i, j)} \cdot (1 - p)^{1 - A_{\mathcal{G}_k}(i, j)}]^{1 - A_{\mathcal{G}^*}(i, j)}. \end{aligned}$$

Hence, the conditional probability of observing a network population $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$ given $A_{\mathcal{G}^*}, p, q$ is,

$$P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | p, q, A_{\mathcal{G}^*}) = \prod_{k=1}^N \prod_{(i,j):i<j} P(A_{\mathcal{G}_k}(i,j) | A_{\mathcal{G}^*}(i,j), p, q) =$$

$$\prod_{k=1}^N \prod_{(i,j):i<j} ((1-q)^{A_{\mathcal{G}_k}(i,j)} \cdot q^{1-A_{\mathcal{G}_k}(i,j)})^{A_{\mathcal{G}^*}(i,j)} \cdot (p^{A_{\mathcal{G}_k}(i,j)} \cdot (1-p)^{1-A_{\mathcal{G}_k}(i,j)})^{1-A_{\mathcal{G}^*}(i,j)}.$$

An advantage of the measurement error formulation is that the model specification for the representative network $A_{\mathcal{G}^*}$ can vary depending on the type of information the analyst wants to capture for the data at hand. As previously discussed, [Le et al. \(2018\)](#) assume a SBM for the network representative. For illustration we also assume an SBM structure for the representative $A_{\mathcal{G}^*}$, but note that this can be easily modified, e.g. reduced to a simpler model such as the Erdős-Rényi if supported by the data. The SBM for the representative can be represented hierarchically in the following way,

$$A_{\mathcal{G}^*}(i,j) | \boldsymbol{\theta}, \mathbf{b} \sim \text{Bernoulli}(\theta_{b_i b_j});$$

$$\theta_{kl} \sim \text{Beta}(\epsilon_{kl}, \zeta_{kl});$$

$$\mathbf{b} | \mathbf{w} \sim \text{Multinomial}(\mathbf{w});$$

where w_k represents the probability of a node to belong to block $k \in \{1, \dots, K\}$. For the probability vector $\mathbf{w} = \{w_1, \dots, w_K\}$ we assume a symmetric Dirichlet prior distribution with hyperparameter $\boldsymbol{\chi}$. Common choices for the hyperparameter vector $\boldsymbol{\chi}$ are setting all elements to 0.5 or 1.

Thus the hierarchical structure of the model is,

$$\prod_{k=1}^N \prod_{(i,j):i<j} P(A_{\mathcal{G}_k}(i,j) | A_{\mathcal{G}^*}(i,j), p, q) P(A_{\mathcal{G}^*}(i,j) | \boldsymbol{\theta}, \mathbf{b})$$

where

$$P(A_{\mathcal{G}^*}(i,j) | \boldsymbol{\theta}, \mathbf{b}) = \theta_{b_i b_j}^{A_{\mathcal{G}^*}(i,j)} (1 - \theta_{b_i b_j})^{1 - A_{\mathcal{G}^*}(i,j)}.$$

We further specify a Beta prior distribution for both the false positive p and false negative q probabilities,

$$p \sim \text{Beta}(\alpha_0, \beta_0), \quad q \sim \text{Beta}(\gamma_0, \delta_0),$$

which facilitates posterior computations. A common choice sets the Beta prior hyperparameters to 0.5, corresponding to the Jeffreys prior.

4.2 Mixture of measurement error models

We further extend the measurement error model to a mixture of measurement error models, with a predefined number of mixture components, C , in order to provide a model-based approach for identifying clusters of networks in a network population $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$. We assume each cluster c of networks is described by a unique network representative $A_{\mathcal{G}_c^*}$, a false positive probability p_c , and a false negative probability q_c , where $c \in \{1, \dots, C\}$. In this section, we present the Bayesian framework for this mixture of measurement error models. Each cluster-specific representative network is characterised by an SBM, and the block structure of each representative is allowed to vary.

Let $\mathbf{z} = (z_1, \dots, z_N) \in \{1, \dots, C\}$ be the latent variables representing the cluster membership of the network data $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$. Then the conditional probability of the data given \mathbf{z} takes the form

$$P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | \{p_c, q_c, A_{\mathcal{G}_c^*}\}_{c=1}^C, z_1, \dots, z_N) =$$

$$\prod_{k=1}^N \left(\prod_{(i,j):i<j} \left((1 - q_{z_k})^{A_{\mathcal{G}_k}(i,j)} q_{z_k}^{1 - A_{\mathcal{G}_k}(i,j)} \right)^{A_{\mathcal{G}_{z_k}^*}(i,j)} \cdot \left(p_{z_k}^{A_{\mathcal{G}_k}(i,j)} (1 - p_{z_k})^{1 - A_{\mathcal{G}_k}(i,j)} \right)^{1 - A_{\mathcal{G}_{z_k}^*}(i,j)} \right).$$

We assume that the cluster labels z_1, \dots, z_N follow a Multinomial distribution with parameter $\boldsymbol{\tau} = (\tau_1, \dots, \tau_C)$, where τ_c represents the probability that a network observation belongs to cluster c , and $\sum_{c=1}^C \tau_c = 1$. We assume a symmetric Dirichlet prior distribution for the vector of probabilities $\boldsymbol{\tau}$ which has the advantage of being conditionally conjugate with the distribution for \mathbf{z} . As commented previously common choices set the Dirichlet hyperparameters all to 0.5 or to 1.

4.3 MCMC scheme for mixture model

With the modelling framework described above we are able to draw samples from the joint posterior distribution of the parameters using MCMC. The joint posterior distribution is known up to a normalising constant, specifically

$$\begin{aligned} & P(\mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z}, \boldsymbol{\tau}, \mathbf{w}, \mathbf{b}, \boldsymbol{\theta} | A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) \\ & \propto P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z}) \cdot P(\mathbf{A}_{\mathcal{G}^*} | \mathbf{w}, \mathbf{b}, \boldsymbol{\theta}) \cdot P(\mathbf{p} | \boldsymbol{\alpha}_0, \boldsymbol{\beta}_0) \\ & \quad \cdot P(\mathbf{q} | \boldsymbol{\gamma}_0, \boldsymbol{\delta}_0) \cdot P(\mathbf{z} | \boldsymbol{\tau}) \cdot P(\boldsymbol{\tau} | \boldsymbol{\psi}) \cdot P(\boldsymbol{\theta} | \boldsymbol{\epsilon}, \boldsymbol{\zeta}) \cdot P(\mathbf{b} | \mathbf{w}) \cdot P(\mathbf{w} | \boldsymbol{\chi}), \end{aligned} \quad (2)$$

where $P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z})$ is the conditional probability of the network data, $P(\mathbf{A}_{\mathcal{G}^*} | \mathbf{w}, \mathbf{b}, \boldsymbol{\theta})$ is the conditional probability of the latent variable $\mathbf{A}_{\mathcal{G}^*}$ and the rest of the components of the right hand side of the expression are the prior distributions for the model parameters, as defined in Sections 4.1 and 4.2.

To obtain posterior samples from the joint posterior in (2) we note that many of the full conditional distributions of the unknown quantities (parameters/latent data) are available in closed form, and when these are not available these can be approximated using Metropolis-Hastings. As a result we obtain posterior inferences through a component wise MCMC sampler, also known as a Metropolis-Hastings-within-Gibbs sampler. This closely follows a Gibbs sampler, where all parameters and latent data are updated from their full conditional distributions except for the full conditionals of $\{A_{\mathcal{G}_c^*}, p_c, q_c\}_{c=1}^C$, which are approximated using Metropolis-Hastings proposal distributions.

In the Metropolis-Hastings step, we use a mixture of kernels for updating the parameters of the measurement error model $\{A_{\mathcal{G}_c^*}, p_c, q_c\}_{c=1}^C$, in analogy to the MCMC scheme seen in Lunagómez et al. (2021). Specifically in every iteration of the MCMC we update the adjacency matrix of the network representative of cluster c , $A_{\mathcal{G}_c^*}$, using either of the following two proposals with some fixed probability:

- (I) We perturb the edges of the current network representative $A_{\mathcal{G}_c^*}^{(curr)}$ of cluster c in the following way:

$$A_{\mathcal{G}_c^*}^{(prop)}(i, j) = \begin{cases} 1 - A_{\mathcal{G}_c^*}^{(curr)}(i, j), & \text{with probability } \omega \\ A_{\mathcal{G}_c^*}^{(curr)}(i, j), & \text{with probability } 1 - \omega \end{cases}.$$

- (II) We propose a new network representative $A_{\mathcal{G}_c^*}^{(prop)}$ for cluster c drawing each edge of the proposed representative $A_{\mathcal{G}_c^*}^{(prop)}(i, j)$ independently from a Bernoulli distribution with parameter $\frac{1}{N} \sum_{k=1}^N A_{\mathcal{G}_k}(i, j)$.

Thus we accept the proposed network representative $A_{\mathcal{G}_c^*}^{(prop)}$ with probability

$$\min \left\{ 1, \frac{P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | A_{\mathcal{G}_c^*}^{(prop)}, p_c^{(curr)}, q_c^{(curr)}, \mathbf{z}^{(curr)}) P(A_{\mathcal{G}_c^*}^{(prop)} | \mathbf{b}_c^{(curr)}, \boldsymbol{\theta}_c^{(curr)})}{P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | A_{\mathcal{G}_c^*}^{(curr)}, p_c^{(curr)}, q_c^{(curr)}, \mathbf{z}^{(curr)}) P(A_{\mathcal{G}_c^*}^{(curr)} | \mathbf{b}_c^{(curr)}, \boldsymbol{\theta}_c^{(curr)})} \cdot \frac{Q(A_{\mathcal{G}_c^*}^{(curr)} | A_{\mathcal{G}_c^*}^{(prop)})}{Q(A_{\mathcal{G}_c^*}^{(prop)} | A_{\mathcal{G}_c^*}^{(curr)})} \right\}, \quad (3)$$

where $P(A_{\mathcal{G}_c^*}^{(\cdot)} | \mathbf{b}_c, \boldsymbol{\theta}_c)$ is the SBM assumed for the representative defined in Section 4.1 and $Q(A_{\mathcal{G}_c^*}^{(\cdot)} | A_{\mathcal{G}_c^*}^{(\cdot)})$ corresponds to the proposal distribution. The proposal distribution under case (I) proposal is symmetric, and so it cancels out from the Metropolis ratio in expression (3).

To update the false positive probability p_c of cluster c , we use a mixture of random walk proposals indexed by l following Lunagómez et al. (2021).

- Draw $v \sim \text{Unif}(-u_l, u_l)$, for $0 < u_l < 0.5$.
- Calculate the candidate proposal value $y = p_c^{(curr)} + v$.
- Propose a new value for p_c (constrained to lie in the interval (0,0.5) for identifiability reasons) as follows,

$$p_c^{(prop)} = \begin{cases} y, & \text{if } 0 < y < 0.5; \\ -y, & \text{if } y < 0; \\ 1 - y, & \text{if } y > 0.5. \end{cases}$$

The mixture is over $\{u_1, \dots, u_L\}$. Thus, we perturb the current state of the false positive probability $p_c^{(curr)}$ using various sizes of u_l , each imposing a less or more drastic change on $p_c^{(curr)}$. We accept the proposed value $p_c^{(prop)}$ with probability

$$\min \left\{ 1, \frac{P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | A_{\mathcal{G}_c^*}^{(curr)}, p_c^{(prop)}, q_c^{(curr)}, \mathbf{z}^{(curr)}) P(p_c^{(prop)} | \alpha_{0,c}, \beta_{0,c})}{P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | A_{\mathcal{G}_c^*}^{(curr)}, p_c^{(curr)}, q_c^{(curr)}, \mathbf{z}^{(curr)}) P(p_c^{(curr)} | \alpha_{0,c}, \beta_{0,c})} \right\}, \quad (4)$$

where $P(p_c^{(\cdot)} | \alpha_{0,c}, \beta_{0,c})$ is a Beta($\alpha_{0,c}, \beta_{0,c}$) prior as in Section 4.1. The proposal distribution for p_c is symmetric, thus it does not appear in the Metropolis ratio in expression (4). In exactly the same manner, we update the false negative probability q_c , for $c \in \{1, \dots, C\}$.

The rest of the parameters are updated via Gibbs samplers, by drawing values from their full conditional posteriors. The full conditional posterior for $\boldsymbol{\tau}$ is given by

$$P(\boldsymbol{\tau} | \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z}, \mathbf{w}, \mathbf{b}, \boldsymbol{\theta}, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) = \text{Dirichlet}(\psi + \eta_1, \dots, \psi + \eta_C). \quad (5)$$

where $\eta_c = \sum_{j=1}^N 1_c(z_j)$, $c = 1, \dots, C$, denotes the number of networks that belong to cluster c .

We draw the latent cluster-membership z_k for each network observation k from a Multinomial distribution with unnormalised probabilities specified in the following way:

$$\begin{aligned} P(z_k = c | \boldsymbol{\tau}, \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, A_{\mathcal{G}_k}) &\propto P(A_{\mathcal{G}_k} | z_k = c, \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}) \cdot P(z_k = c | \boldsymbol{\tau}) \\ &= \tau_c \cdot \prod_{(i,j): i < j} \left((1 - q_c)^{A_{\mathcal{G}_k}(i,j)} q_c^{1 - A_{\mathcal{G}_k}(i,j)} \right)^{A_{\mathcal{G}_c^*}(i,j)} \cdot \left(p_c^{A_{\mathcal{G}_k}(i,j)} (1 - p_c)^{1 - A_{\mathcal{G}_k}(i,j)} \right)^{1 - A_{\mathcal{G}_c^*}(i,j)} \end{aligned} \quad (6)$$

where $P(A_{\mathcal{G}_k} | z_k = c, \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q})$ is the probability we observe network k given cluster membership $z_k = c$, described by a measurement error model. The normalised probabilities are obtained via Bayes Theorem.

The full conditional posterior for \mathbf{w}_c is

$$P(\mathbf{w}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \boldsymbol{\theta}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) = \text{Dirichlet}(\chi + h_1, \dots, \chi + h_K).$$

where h_k denotes the number of the nodes that belong to block k .

The full conditional posterior for the vector of the block-specific probabilities of an edge occurrence for the network representative of cluster c , $\boldsymbol{\theta}_c$, is

$$P(\boldsymbol{\theta}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \mathbf{w}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) = \text{Beta}(A_{\mathcal{G}_c^*}[kl] + \epsilon_{kl}, \zeta_{kl} + n_{c,kl} - A_{\mathcal{G}_c^*}[kl]). \quad (7)$$

where $A_{\mathcal{G}_c^*}[kl] = \sum_{(i,j): b_{c,i}=k, b_{c,j}=l} A_{\mathcal{G}_c^*}(i,j)$ represents the sum of the entries for the pairs of nodes of the network representative for cluster c that have block membership k, l respectively, and $n_{c,kl} = \sum_{(i,j): i \neq j} \mathbb{I}(b_{c,i} = k, b_{c,j} = l)$ represents the number of the pair of nodes of the representative for cluster c that have membership k, l respectively.

Similarly to the formulation obtained for updating the latent cluster-membership \mathbf{z} of the network data, we obtain updates of the latent block-membership \mathbf{b}_c for the nodes of the network representative of cluster c from a Multinomial distribution with unnormalised probabilities specified as follows:

$$\begin{aligned} P(b_{c,i} = k | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \boldsymbol{\theta}_c, \mathbf{w}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) &\propto P(A_{\mathcal{G}_c^*} | \mathbf{w}_c, \boldsymbol{\theta}_c, b_{c,i} = k) \cdot P(b_{c,i} = k | \mathbf{w}_c) \\ &= w_{c,k} \cdot \prod_{j=1}^n \theta_{kb_{c,j}}^{A_{\mathcal{G}_c^*}(i,j)} (1 - \theta_{kb_{c,j}})^{1 - A_{\mathcal{G}_c^*}(i,j)}. \end{aligned} \quad (8)$$

where $P(A_{\mathcal{G}_c^*} | \mathbf{w}_c, \boldsymbol{\theta}_c, b_{c,i} = k)$ is the probability of observing the representative of cluster c , $A_{\mathcal{G}_c^*}$, described by an SBM, given its i^{th} node belongs to block k . Normalised probabilities are obtained by Bayes Theorem.

For the detailed derivation of the full conditional posterior distributions refer to the Supplementary material Section 1 (Mantziou et al., 2023). In addition, the MCMC algorithm for clustering is sketched in the Supplementary material Algorithm 1 (Mantziou et al., 2023).

4.4 Outlier network detection

Motivated by the application on brain networks, in this section we present a modification of the mixture model presented in Section 4.2 to further explore the heterogeneity in a network population. Specifically, we modify our mixture model to detect a cluster of outlier networks that are different to the majority of the networks in the population. Under this formulation, we are able to address additional applied research questions of interest. In particular, we would like to be able to identify individuals with different brain connectivity patterns compared to the rest of the population.

In contrast to the mixture model formulated for multiple cluster representatives, the outlier cluster detection model assumes a single network representative for the whole population of networks. Under this setup, we assume that there are ultimately two clusters of networks formed within the population of networks, one cluster being the majority cluster, while the other cluster determining the outlier networks in the population. Thus, while the false positive and false negative probabilities remain component specific for each of the two clusters, the network representative is no longer a component specific latent variable.

Similar to the mixture model formulated in Section 4.2, we now specify the number of clusters to $C = 2$, and $\mathbf{z} = (z_1, \dots, z_N) \in \{1, 2\}$ denotes the latent cluster membership of the network data $A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}$. Under the assumption of a single network representative, $A_{\mathcal{G}^*}$, the conditional probability of the data given the latent variables, \mathbf{z} , takes the form,

$$P(A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N} | \{p_c, q_c\}_{c=1}^C, A_{\mathcal{G}^*}, z_1, \dots, z_N) = \prod_{k=1}^N \left(\prod_{(i,j):i<j} \left((1 - q_{z_k})^{A_{\mathcal{G}_k(i,j)}} q_{z_k}^{1 - A_{\mathcal{G}_k(i,j)}} \right)^{A_{\mathcal{G}^*(i,j)}} \cdot \left(p_{z_k}^{A_{\mathcal{G}_k(i,j)}} (1 - p_{z_k})^{1 - A_{\mathcal{G}_k(i,j)}} \right)^{1 - A_{\mathcal{G}^*(i,j)}} \right).$$

Again, the model specification for the representative network can vary depending on the type of information we want to capture for the data at hand. A common choice is to consider an SBM structure again for the representative. Due to having now a single representative only, the SBM model parameters are no more component specific to the cluster.

To sample from the joint posterior of this model, we develop a Metropolis-Hastings-within-Gibbs MCMC scheme, as presented in Section 4.3. The full conditional posterior distributions are obtained as seen in Section 4.3 with the only difference that the parameters/latent variables characterising the representative, $A_{\mathcal{G}^*}$, namely, the block-specific edge probabilities, $\boldsymbol{\theta}$, the probability of a node to belong to a block, \mathbf{w} , and the block membership of the nodes, \mathbf{b} , are no longer component specific, i.e. not indexed by cluster c .

4.5 Sparse Finite Mixture extension

In practice, the number of clusters is not known a priori, and so we need to be able to determine an appropriate number of clusters to specify in our model. This is a problem that has been extensively discussed in the literature, and a variety of approaches have been proposed. For finite mixture models, a common approach for estimating C is through information criteria such as BIC (Keribin, 2000). Alternative approaches can incorporate uncertainty in the number of clusters such as reversible jump MCMC methods (Richardson and Green, 1997) or Dirichlet Process (DP) mixture models (Neal, 2000), with the former being particularly challenging in network models, due to the MCMC moving between very different dimensions.

We adopt the approach in Malsiner-Walli et al. (2016), who develop a method that conveniently extends the finite mixture model to make inferences with an unknown number of clusters, known as the Sparse Finite Mixture (SFM) model. A sparse symmetric Dirichlet prior distribution is specified for the weights of the mixture components, and the number of mixture components are assumed to be more than the number of clusters in the data. Frühwirth-Schnatter and Malsiner-Walli (2019) discuss how the SFM model can be easily extended to examples with non-Gaussian data, and make comparisons to DP mixtures with respect to their clustering performance. The convenient implementation of the SFM model for finite mixture models, together with its wide applicability for different types of data, led us to consider this as an extension to our finite mixture model that allows an unknown number of clusters C .

Extending our finite mixture model to the SFM model requires the specification of a symmetric Dirichlet prior distribution $\text{Dir}(e_0, \dots, e_0)$ of order C_{max} , with C_{max} being an upper bound on the number of clusters such that $C < C_{max}$ where C is the number of clusters in the data, resulting in an overfitted mixture model. The size of e_0 should result in many of the weights $\boldsymbol{\tau} = (\tau_1, \dots, \tau_{C_{max}})$ being close to 0, imposing sparsity on the number of clusters. We specify a Gamma(a_e, b_e) hyperprior on the hyperparameter, e_0 , of the Dirichlet prior, and sample from the posterior of e_0 using a Metropolis-Hastings step as proposed

in [Frühwirth-Schnatter and Malsiner-Walli \(2019\)](#). Specifically, in our model we specify a Random Walk proposal for the MH step for e_0 . The specification of a_e, b_e plays a key role in the clustering performance of the model and should result in values of e_0 close to 0.

In [Section 5.3](#), we explore the performance of the SFM extension of our model on simulated network populations, as well as implement it in the real data applications in [Section 6](#).

5 Simulations

In this Section, we perform simulation studies to assess the performance of our algorithm in inferring the model parameters/latent variables and clustering network data. First, we explore the performance of our algorithm for moderate network sizes and various noise levels and SBM models, and second, we investigate the algorithm performance for various network and sample sizes.

5.1 Moderate network sizes

In this simulation study we investigate the performance of our model in inferring model parameters for network populations with a moderate number of nodes in different scenarios. Specifically, we consider the case of networks with $n = 21$ nodes, and a population of $N = 180$ networks. We assume $C = 3$ clusters of networks in the population, and consider SBMs for each representative network with $B = 2$ blocks. We vary the model parameter values in order to explore performance.

To simulate the network population, we first simulate the representative network of each cluster. We generate representatives of the clusters under two different SBM structures with parameters as shown in [Table 1](#) (left). In [Figure 3](#) we visualise the 21-node representatives for each of the $C = 3$ clusters under each of the SBM structures 1 and 2 assumed.

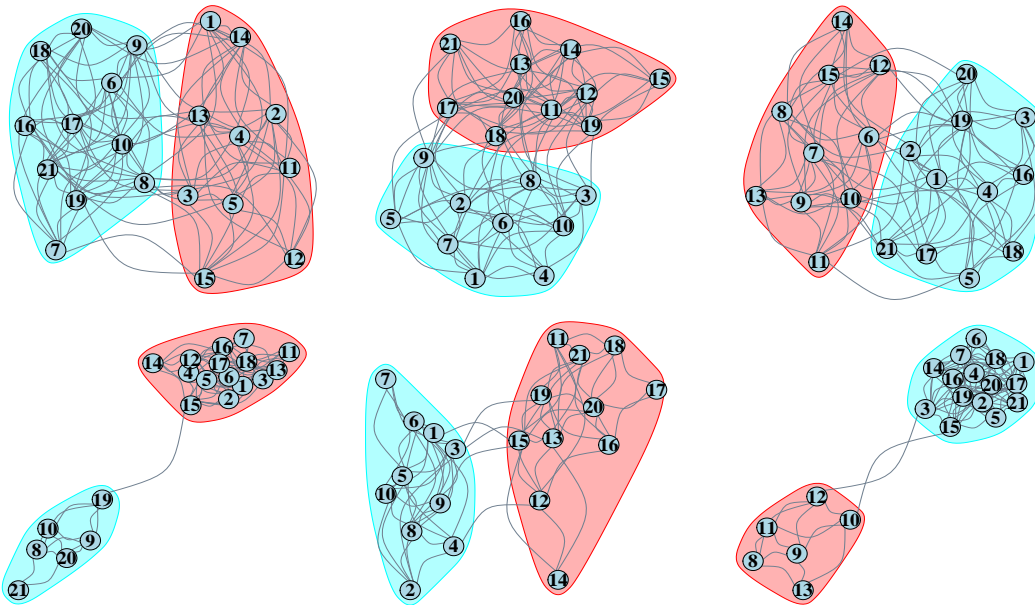


Figure 3: Top: Network representatives for clusters $c = 1, 2$ and 3 respectively (from left to right), under SBM structure 1. Bottom: Network representatives for clusters $c = 1, 2$ and 3 respectively (from left to right), under SBM structure 2.

Next, we generate a population of 180 networks by perturbing the edges of each representative through a measurement error process. Specifically, we generate edges for 60 networks in each of the $C = 3$ clusters, depending on the existence or non-existence of an edge in the representative network of the corresponding cluster c , given a false positive p_c and false negative q_c probability. The simulation regimes considered for p_c, q_c are presented in [Table 1](#) (right).

For each simulation regime, we run our MCMC for 500,000 iterations with a burn-in of 150,000. In [Figure 4](#), we visualise the posterior distribution of p_c and q_c under the simulation regime with $p_c = 0.1$ and $q_c = 0.3$

SBM	c	θ			w		SBM _{i}	
		θ_{11}	θ_{12}	θ_{22}	w_1	w_2	p_c	q_c
1	1	0.8	0.2	0.8	0.5	0.5	0.1	0.2
	2	0.8	0.2	0.8	0.5	0.5	0.1	0.3
	3	0.8	0.2	0.8	0.5	0.5	0.2	0.1
2	1	0.7	0.05	0.8	0.7	0.3	0.2	0.3
	2	0.7	0.05	0.8	0.5	0.5	0.3	0.1
	3	0.7	0.05	0.8	0.3	0.7	0.3	0.2

Table 1: Simulation regimes for 21-node networks and $C = 3$ clusters. Left Table: SBM structures 1 and 2 for simulating a network representative for each cluster c . Right Table: sizes of false positive p_c and false negative q_c probabilities used to simulate network data under each SBM structure.

for all c , and SBM structure 2 for the representative networks. In addition, in Figure 5, we visualise the posterior distribution of p_c and q_c under the simulation regime with $p_c = 0.2$ and $q_c = 0.3$ for all c , and SBM structure 1 for the representative networks. The bar in the violin plot indicates the 95% credible interval and the point indicates the posterior mean. In both Figures, we observe that the posterior means are very close to the true values of the parameters, while the 95% credible intervals enclose the true values of the parameter for all cases. This finding also holds for the rest of the simulation regimes. In Supplementary Material, Section 2.1 (Mantziou et al., 2023), we summarise the results obtained for the rest of the simulation regimes in Tables 1-15.

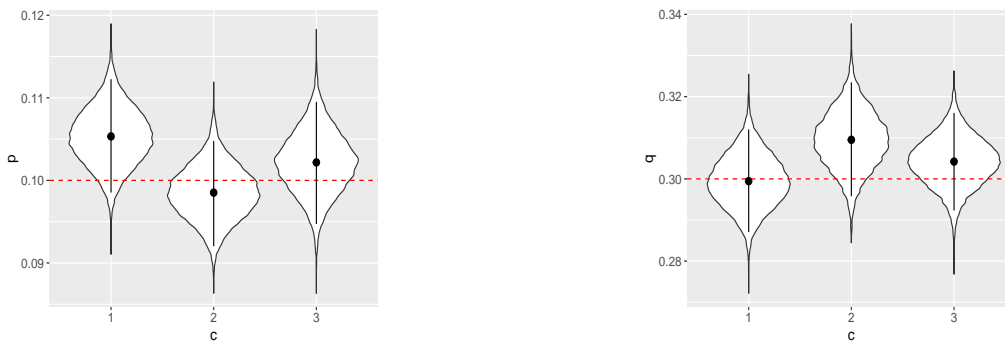


Figure 4: Posterior distribution of false positive probabilities p_c (left) and false negative q_c (right) for $c \in \{1, 2, 3\}$, for simulation regime with $p_c = 0.1$ and $q_c = 0.3$. Red dotted lines indicate the true value of p_c and q_c .

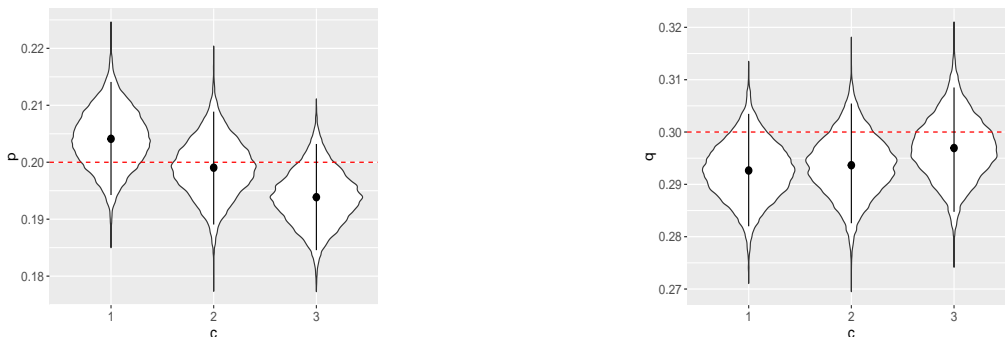


Figure 5: Posterior distribution of false positive probabilities p_c (left) and false negative q_c (right) for $c \in \{1, 2, 3\}$, for simulation regime with $p_c = 0.2$ and $q_c = 0.3$. Red dotted lines indicate the true value of p_c and q_c .

In order to investigate the performance of our algorithm in identifying the true representatives, we obtain the Hamming distance between the posterior representative samples and the true representatives, after a

burn-in of 150,000 and a lag of 50, leaving 7,000 posterior samples. The Hamming distance measures how dissimilar two graphs are with respect to their edges (Donnat and Holmes, 2018). The maximum Hamming distance between two undirected, 21-node networks is equal to $\binom{21}{2} = 210$, meaning the two graphs have no edges in common. We calculate the proportion of times that the distance is less or equal to 1, 5 and 10 respectively, similar to the summaries obtained in Lunagómez et al. (2021). For each simulation regime, we observe that all the posterior representative samples drawn for each cluster have a Hamming distance from the true representative less than or equal to 1, 5, and 10, 100% of the time, as presented in the Supplementary material, Section 2.1, Tables 16-17 (Mantziou et al., 2023). This result suggests that the true representatives are almost fully identified from our algorithm.

In addition, we assess the effectiveness of our algorithm in identifying the cluster membership of the networks using the clustering entropy and purity indices. Both clustering entropy and clustering purity are indices for evaluating clustering performance when the true cluster labels are known (Kim and Park, 2007; Schütze et al., 2008). We note that a clustering entropy value of 0 and clustering purity value of 1 indicate a perfect cluster allocation of the networks. We obtain 7,000 posterior draws (after a burn-in of 150,000 and lag of 50) for the cluster membership z , calculate the clustering entropy and clustering purity with respect to the true membership of the networks, and calculate their mean for each simulation regime. The simulation results indicate the true cluster membership of the networks is fully recovered by our MCMC algorithm, with mean entropy 0 and mean purity 1 for each simulation regime. These results are in the Supplementary material, Section 2.1, Table 18 (Mantziou et al., 2023).

We further compare our approach to the nonparametric Bayesian approach for modelling network populations by Durante et al. (2017), and to the maximum likelihood approach for clustering network populations by Signorelli and Wit (2020). Durante et al. (2017) were originally interested in flexibly modelling a population of networks with diverse characteristics, rather than clustering network-valued data, although clustering is a natural extension of their approach. Both Durante et al. (2017) and Signorelli and Wit (2020) capture the heterogeneity in a network population through a model-based framework, as is also the case in our framework. However, neither approach permits easy interpretation of the clusters. The fundamental difference in our approach is the ability to infer a network representing each cluster, thus providing a useful summary of the networks in each cluster which can be advantageous when making inferences for diverse real-data applications.

We implement both approaches on the diverse network populations simulated as described earlier in this section, and assess the performance of the methods in clustering the network observations with respect to the clustering entropy and purity indices. Durante et al. (2017) method perfectly recovers the underlying $C = 3$ clusters in the simulated network populations, with mean clustering entropy 0 and mean clustering purity 1 for each simulation regime.

To implement Signorelli and Wit (2020) we need to first specify a statistical network model for the components of their mixture model. We choose to use an SBM as it is the model that we assumed for the network representatives of the clusters for simulating the network populations. However, there are two key limitations in the approach of Signorelli and Wit (2020). First, it is assumed that all clusters share the same block structure, and second, the block structure should be pre-specified as it is not inferred in their scheme, in contrast to our model which infers the block structure of the representatives and allows it to vary between the network representatives. Thus, to obtain a single block structure for all three mixture components, we use majority vote to determine the block membership of each node using the block structures specified for the representatives in our simulations, for each SBM simulation scenario 1 and 2. The clustering entropy and clustering purity calculated from the results obtained after applying the mixture model of Signorelli and Wit (2020), for each simulation regime, are presented in Table 2. We see our model outperforms Signorelli and Wit (2020) on our simulated data, which can be attributed to their model’s restrictive assumption of a single SBM structure for all mixture components. As a result we do not consider the approach of Signorelli and Wit (2020) any further.

We now explore the performance of our model on data simulated under a different parameterisation to a mixture of measurement error models. Specifically, we consider the simulated population of networks in Durante et al. (2017). The data comprises 100 networks, with each network generated using one of four possible parameterisations of the edge probabilities, resulting in four groups of networks with different properties. Specifically, the four groups are characterised by a community structure, small-worldness, an Erdős-Rényi structure and scale-free properties respectively. We run our MCMC for 500,000 iterations and assess the accuracy of our algorithm in inferring the group membership of each network in the population using the clustering entropy and clustering purity indices described earlier. Specifically, we consider 7,000

		SBM ₁		SBM ₂	
p_c	q_c	Entropy	Purity	Entropy	Purity
0.1	0.2	0.69	0.7	0.71	0.69
	0.3	0.79	0.57	0.61	0.78
0.2	0.1	0.65	0.73	0.55	0.81
	0.3	0.79	0.57	0.41	0.86
0.3	0.1	0.76	0.63	0.56	0.81
	0.2	0.75	0.65	0.59	0.74

Table 2: Clustering entropy and purity for results obtained after implementing [Signorelli and Wit \(2020\)](#) model on the simulated data presented in this section.

posterior draws (after a burn-in of 150,000 iterations and lag of 50) for the cluster membership z of the networks, and obtain mean clustering entropy of 0 and mean clustering purity of 1. This indicates our algorithm perfectly recovers the group membership of the networks in the population, despite the different generative process used to simulate the data.

The simulation results so far show perfect performance of our model in recovering the true cluster labels for the network observations. To investigate clustering performance under more challenging scenarios, we consider simulating network populations with high false positive and negative probabilities ($p_c = q_c = 0.4$) for a range of smaller population sizes, ranging from 36 to 180. In each population, $C = 3$ and network representatives have SBM structure 1 (Figure 3 top). We compare results to those obtained using the model in [Durante et al. \(2017\)](#) on the same data. For our model, the MCMC is run for 500,000 iterations with 150,000 burn-in and lag 87 resulting in 4,023 MCMC draws, while for the [Durante et al. \(2017\)](#) approach the MCMC is run for 5,000 iterations with 1,000 burn-in, leaving 4,000 posterior draws. Figures 6 and 7 illustrate the distribution of the clustering entropy and clustering purity for the posterior draws under each model, as well as for different sample sizes.

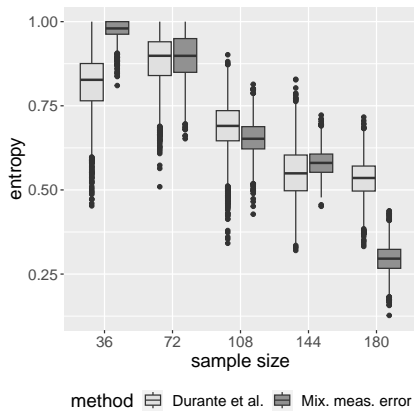


Figure 6: Distribution of the clustering entropy across posterior draws, for our method and [Durante et al. \(2017\)](#), for varying sample sizes.

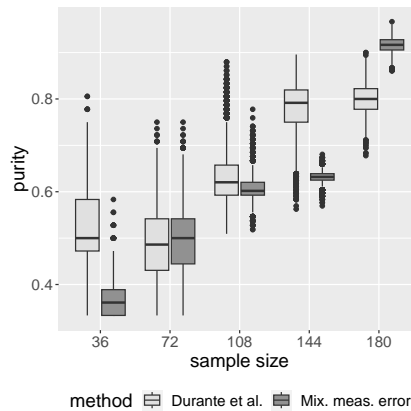


Figure 7: Distribution of the clustering purity across posterior draws, for our method and [Durante et al. \(2017\)](#), for varying sample sizes.

We observe that the clustering performance of the two methods have similar relationships with the sample size. As expected, both approaches perform least well in recovering the true cluster configurations for high noise levels and the smallest sample sizes, but steadily improve as sample sizes increase. We observe a slightly better performance with [Durante et al. \(2017\)](#) over our approach for the smallest sample size of 36 networks, with the converse for the biggest sample size of 180 networks. It is worth noting again that although both approaches have similar clustering performance overall, a key advantage of our method is the interpretability of the clusters with respect to a network representative, which is particularly relevant for our motivating data applications.

We additionally investigate the performance of our model in the network population size of 180 for various different noise levels. Specifically, we generate network populations by perturbing the edges of the representatives of SBM structure 1, illustrated in Figure 3, with varying sizes of the false positive and negative

probabilities given in Tables 3 and 4.

p_c	q_c
0.01	
0.05	
0.1	
0.15	
0.2	0.1
0.25	
0.3	
0.35	
0.4	
0.45	

Table 3: Simulation regimes for varying sizes of false positive probabilities p_c and fixed false negative probabilities q_c , for $c \in \{1, 2, 3\}$ and 21-node networks.

p_c	q_c
	0.01
	0.05
	0.1
	0.15
0.1	0.2
	0.25
	0.3
	0.35
	0.4
	0.45

Table 4: Simulation regimes for varying sizes of false negative probabilities q_c and fixed false positive probabilities p_c , for $c \in \{1, 2, 3\}$ and 21-node networks.

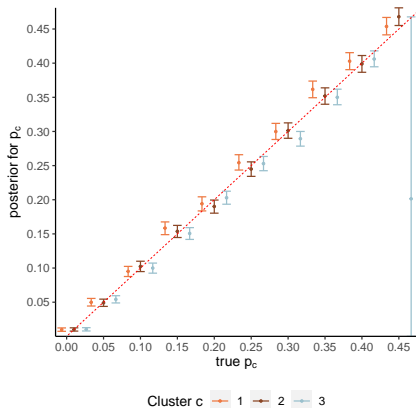


Figure 8: Posterior means and 95% credible intervals for false positive probabilities p_c for $c \in \{1, 2, 3\}$ (y axis), plotted against the true values of p_c (x axis). Red dashed line is the $y=x$ line.

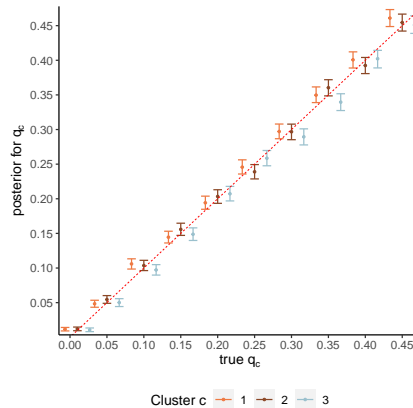


Figure 9: Posterior means and 95% credible intervals for false negative probabilities q_c for $c \in \{1, 2, 3\}$ (y axis), plotted against the true values of q_c (x axis). Red dashed line is the $y=x$ line.

For each regime presented in Tables 3 and 4, we generate a network population and run the MCMC for 500,000 iterations with a burn-in of 150,000 iterations. In Figures 8 and 9, for each cluster $c \in \{1, 2, 3\}$, we plot the posterior means and 95% credible intervals (via errors bars) for the different false positive (Table 3) and false negative probabilities (Table 4) respectively against their true values. We see the posterior means lie mostly close to the $y = x$ line (red dashed line), indicating our model performs well in inferring the true false positive and false negative probabilities, even for high noise levels. However, in Figure 8, for the highest noise value of $p_c = 0.45$ for $c \in \{1, 2, 3\}$, we observe that the posterior mean of the false positive probability of cluster 3, p_3 , is equal to 0.2, which is substantially different to its true value, while its 95% credible interval covers a wide range of values, indicating that our MCMC chain struggles to make inferences here. These results suggest that our model performs well in most cases, even for high noise levels, but we must be cautious when making inferences for network populations with great variability in their structure.

For the simulations reported in this section, the computational time required to run our MCMC procedure for 500,000 iterations varied from approximately 50 minutes (for a population size of 36 networks) through to approximately 80 minutes (for a population size of 180 networks). We note here that in both our simulations and data analysis, we consider a large number of iterations to ensure convergence of our MCMC, even though for some scenarios less iterations would suffice.

5.2 Varying sizes of networks and network populations

We now explore how well our model infers the parameters with respect to various network sizes and sample sizes. We keep $C = 3$ clusters and $B = 2$ blocks. We consider four different network sizes of 25, 50, 75 and 100 nodes, and simulate populations of 45, 90, 135, 180, 225, 270 and 315 networks, for each network size respectively.

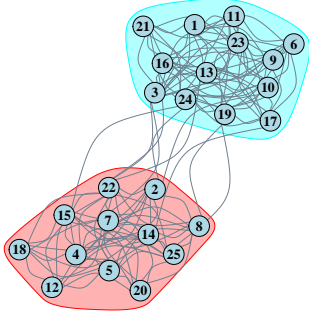


Figure 10: 25-node representative of cluster labelled 1.

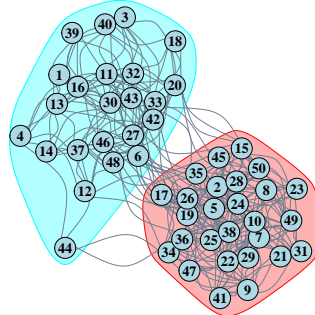


Figure 11: 50-node representative of cluster labelled 1.

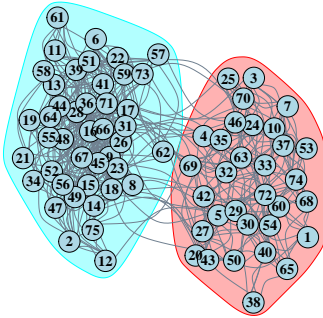


Figure 12: 75-node representative of cluster labelled 1.

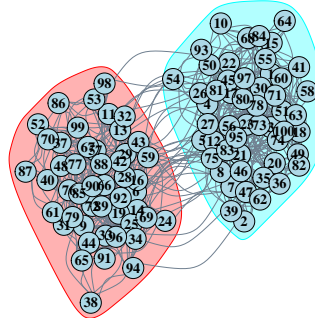


Figure 13: 100-node representative of cluster labelled 1.

To simulate the network populations, we first generate the network representatives of each cluster assuming each follows an SBM. We specify the parameters of the SBMs so that the expected degree of the resulting network representatives is preserved for all network sizes. The representatives obtained for cluster $c = 1$ for the different network sizes are visualised in Figures 10, 11, 12 and 13. In the Supplementary material, Section 2.2, Figures 1-4 (Mantziou et al., 2023), we illustrate the rest of the representatives for $c = 2, 3$, for the network sizes considered. The resulting populations are generated by perturbing the edges of each network representative, for each network size considered, with a false positive p_c and false negative q_c probability fixed at 0.08, for $c \in \{1, 2, 3\}$.

For each simulation regime, we consider 10 replications of our MCMC, each on a different randomly generated data set, and run our algorithm for 500,000 iterations with a burn-in of 150,000. For the largest simulation scenario involving 100-node networks and 315 network observations, 500,000 iterations of our MCMC required a run time of approximately 24 hours. Thus, 10 replications was deemed reasonable, taking into account the increased computational burden associated with running these simulations. We demonstrate the performance of our model by obtaining the distribution of the absolute error of the model parameters for each simulation regime, as seen in Figures 14 and 15. Specifically, the plots demonstrate how the distribution of the absolute error (y axis) scales for various sample sizes (x axis). The absolute error is the absolute value of the difference between the posterior means obtained after burn-in, and the true value of the parameter. The different grey shades of the boxplots correspond to the different clusters considered, and the lines connect the mean absolute error across replications to illustrate the trend.

The plots indicate that the sample size affects the performance of our model for the network sizes considered. We observe that as the number of nodes increase, there is an increase in the required number of

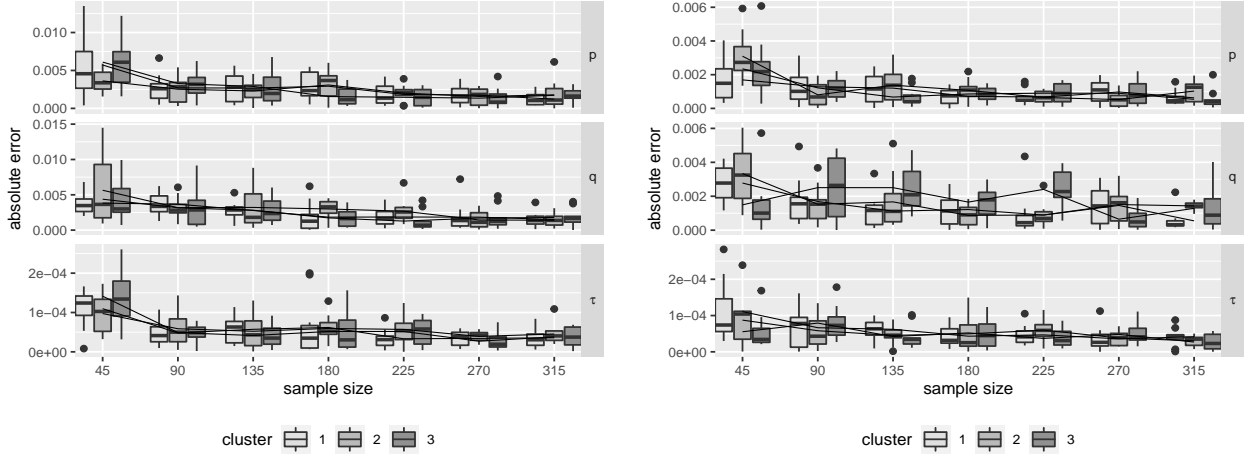


Figure 14: Left: Absolute error (y axis) for model parameters p , q and τ , for 25-node networks and varying population sizes (x axis). Right: Absolute error (y axis) for model parameters p , q and τ , for 50-node networks and varying population sizes (x axis).

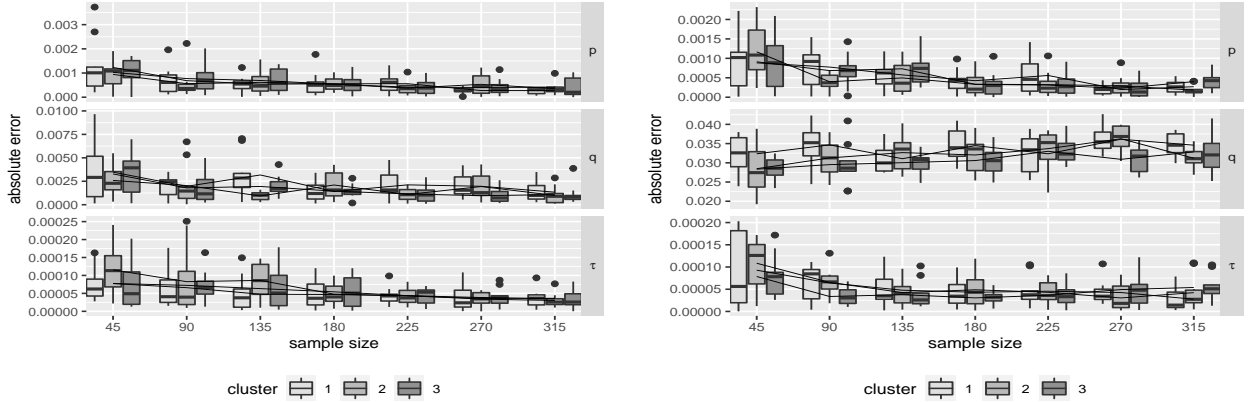


Figure 15: Left: Absolute error (y axis) for model parameters p , q and τ , for 75-node networks and varying population sizes (x axis). Right: Absolute error (y axis) for model parameters p , q and τ , for 100-node networks and varying population sizes (x axis).

networks to preserve the same level of accuracy in estimation. Nevertheless, even for large networks and small sample sizes, it is encouraging to see the posterior means obtained are not far away from their true values. We further illustrate the performance of our model in identifying the true representatives of each cluster for the various sample and network sizes, in the Supplementary Material, Section 2.2, Figure 5 (Mantziou et al., 2023).

For the simulation regime with the largest size of networks and network population, $n = 100$ nodes and $N = 315$ networks respectively, we additionally explore the performance of our model for different hyperparameter settings for the Beta and Dirichlet priors of our model. Specifically, we consider two common non-informative priors, the Jeffrey’s prior, Beta(0.5,0.5) and Dirichlet(0.5), and the Uniform prior, Beta(1,1) and Dirichlet(1). Additionally, we specify an informative prior such that the proportional reduction in variance from the Uniform prior to the informative prior is equal to the proportional reduction in variance from the Jeffrey’s prior to the Uniform prior. This results in Beta(1.75,1.75) and Dir(1.75) priors. Under each hyperparameter setting, we run our MCMC for 500,000 iterations and obtain similar results, not presented herein due to length restrictions, indicating that our model is not sensitive to the hyperparameter specification.

Lastly, for completeness, we explore the clustering performance of the model of Durante et al. (2017) for varying sizes of networks. Specifically, we implement the model by Durante et al. (2017) on the simulation regimes with $N = 180$ and $n = 25, 50, 75, 100$. The model accurately clusters the networks with mean

clustering entropy 0 and mean clustering purity 1 for each regime considered. As previously noted, despite the good clustering performance of this model, meaningful interpretations of the clusters can be challenging.

5.3 Number of clusters with SFM

In the simulations performed in Sections 5.1 and 5.2, the number of clusters were known and specified according to the number of mixture components used to simulate the network populations, which is typically not the case in real data applications. The SFM extension of our model introduced in Section 4.5 allows us to treat the number of clusters as unknown within our framework, and subsequently infer an appropriate number of clusters to specify.

To explore the performance of the SFM extension of our model in identifying the true number of clusters and the true cluster labels of the networks in a population, we implement it on a range of simulation regimes described in Section 5.1. Specifically, we consider the simulated network data under SBM structure 1 and noise levels $p_c = 0.1$ and $q_c = 0.2$ as well as the case with noise levels $p_c = 0.2$ and $q_c = 0.3$, for $c \in \{1, 2, 3\}$. Similarly, we consider the simulated network populations under the same noise levels, for the representatives generated with SBM structure 2. The true number of clusters in all cases are $C = 3$. We tune our MCMC specifying $C_{max} = 10$ clusters and hyperparameters $a_e = 1$ and $b_e = 400$ for the Gamma prior on e_0 to impose strong shrinkage of e_0 to 0, and run our MCMC algorithm for 500,000 iterations.

We notice we quickly converge to the true number of clusters, $C = 3$, not illustrated herein due to length restrictions. To assess the model performance in identifying the true cluster labels of the networks observations, we calculate the clustering entropy and clustering purity indices for \mathbf{z} after a burn-in of 150,000 iterations and a lag of 50, leaving 7,000 posterior draws. We observe that for all four simulation regimes considered in this study, the model perfectly recovers the cluster membership of the networks with mean clustering entropy 0 and mean clustering purity 1.

We note that the SFM model is highly sensitive to the tuning of a_e, b_e in the Gamma hyperprior and C_{max} , as also discussed in Frühwirth-Schnatter and Malsiner-Walli (2019). Different levels of shrinkage can lead to inferences with different numbers of clusters C and cluster configurations. In the next section, involving the motivating real data applications, we consider two ways of determining the number of clusters C , one of which is the SFM model.

6 Motivating data examples

In this Section, we present the application of our mixture model on the two real-world populations of network data sets presented in Section 2.

6.1 Movement patterns across campus

As introduced in Section 2.1, Tacita is a mobile phone application that records the displays visited by users, along with the time visited and the type of content shown on the display. One way to represent the data collected from the Tacita application is through a network, where nodes correspond to displays, and edges correspond to movements of users among the displays. Consequently, we obtain a population of network data set where each network observation corresponds to a user’s movements across displays. The final data sample consists of 120 undirected and unweighted network observations that share the same set of 37 nodes corresponding to the displays across campus. Names of the displays are presented in the Supplementary Material, Section 3.1, Table 19 (Mantziou et al., 2023). As our mixture model requires the pre-specification of the number of clusters C , one way to determine an appropriate number of clusters in the data is through exploratory data analysis (EDA). We considered various network distance metrics, and for each metric we obtain a distance matrix that contains the pairwise distances of the networks in the population. We obtain the Multi Dimensional Scaling (MDS) plot for each distance matrix obtained. The MDS algorithm maps objects in a 2-d space, respecting their pairwise distances. The MDS plots obtained from the EDA are presented in the Supplementary Material, Section 3.1, Figures 6 and 7 (Mantziou et al., 2023).

As anticipated different network distances reveal different type of similarities between the networks. In this regard, the MDS visualisations provide only a point of reference to determine the number of clusters. To begin with, the specification of $C = 3$ seems a reasonable starting point for our analysis as per the EDA

results. Later in this section, we compare this with the SFM extension of our model presented in Section 5.3. To meaningfully initialise the networks’ cluster membership, we combine the results from four different distance metrics; the Hamming, the Jaccard, the l_2 , and the wavelets metrics. For a descriptive review on the distance metrics refer to [Donnat and Holmes \(2018\)](#). Specifically, we use a k-means algorithm using the R package `kmed` [Budiaji \(2019\)](#) to determine four different cluster memberships, corresponding to the four different metrics considered, and determine the final cluster membership initialisation using majority vote, i.e. by determining which networks are consistently allocated to one of the three clusters among the four memberships obtained. We initialise the representative of each cluster by generating its edges using independent Bernoulli draws. The probability with which we draw an edge between two specific nodes of the representative corresponds to the proportion of times that we see that edge in the network data of the corresponding cluster. We then initialise the nodes’ block membership of the representatives using SBM estimates from the R package `blockmodels` [INRA and Leger \(2015\)](#) that suggest the presence of two underlying blocks. We note here that a simpler network model, namely the Erdős-Rényi model, could also be applied to describe the representative networks. We run our MCMC for 500,000 iterations with a burn-in of 100,000.

In Figure 16, the left and middle plots present the proportion of times that the nodes of the representative networks of clusters 2 and 3 respectively, belong in each of the two blocks specified. The results for the representative of cluster $c = 1$ are presented in the Supplementary material, Section 3.1, Figure 10 ([Mantziou et al., 2023](#)). We note here, that a block structure is not identified for the representative of cluster $c = 1$. From Figure 16, we observe a similar block membership is revealed for the representatives of clusters 2 and 3. However, we notice differences in the block allocation of some nodes between the two representatives, namely nodes labelled 1, 3, 5, 7, 10, 11, 12 and 13. In addition, for the representative of cluster $c = 3$, the nodes are more clearly allocated to each of the two blocks, as seen from the proportions. In Figure 16, the right plot corresponds to the proportion of times that an individual is allocated to the clusters, showing that most individuals are clearly allocated to one of the three clusters.

In addition, in Figure 17 we obtain the network representatives of each cluster, with node labels corresponding to the numbering seen in Table 19 given in the Supplementary material, Section 3.1 ([Mantziou et al., 2023](#)), and layout similar to the true location of the displays on campus. The two different colours of the nodes denote the block membership inferred for each representative. In the Supplementary Material, Section 3.1, Figures 8 and 9 ([Mantziou et al., 2023](#)), we present trace plots of the false negative and false positive probabilities for each of the three clusters.

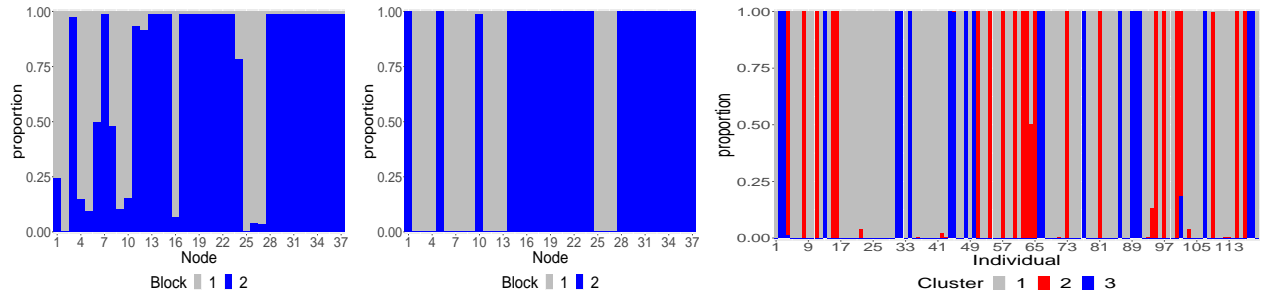


Figure 16: Left: Proportion of times (y axis) that each node (x axis) of the representative network of cluster labelled 2 is allocated in Blocks 1 or 2, after a burn-in of 100,000 iterations. Middle: Proportion of times (y axis) that each node (x axis) of the representative network of cluster labelled 3 is allocated in Blocks 1 or 2, after a burn-in of 100,000 iterations. Right: Proportion of times (y axis) an individual’s network (x axis) is allocated to each of the 3 clusters, after a burn-in of 100,000 iterations.

For cluster $c = 1$, we observe that the representative concentrating the whole posterior mass is sparse having only few edges and no SBM structure. Specifically, we note that the edges of the representative correspond to movements of individuals among displays that are very close to each other (e.g. edge between nodes 4 and 10 corresponding to displays both located in the same building). In addition, most of the networks in the population are allocated to this cluster with a very small false positive probability with posterior mean 0.003 and very large false negative probability with posterior mean 0.49. The small false positive probability indicates that the edges observed in the network data are correctly recorded, while the high false negative probability indicates a high possibility of edges in the network data that we do not observe. This is a reasonable finding as we would anticipate that the Tacita application might have missed movements of users

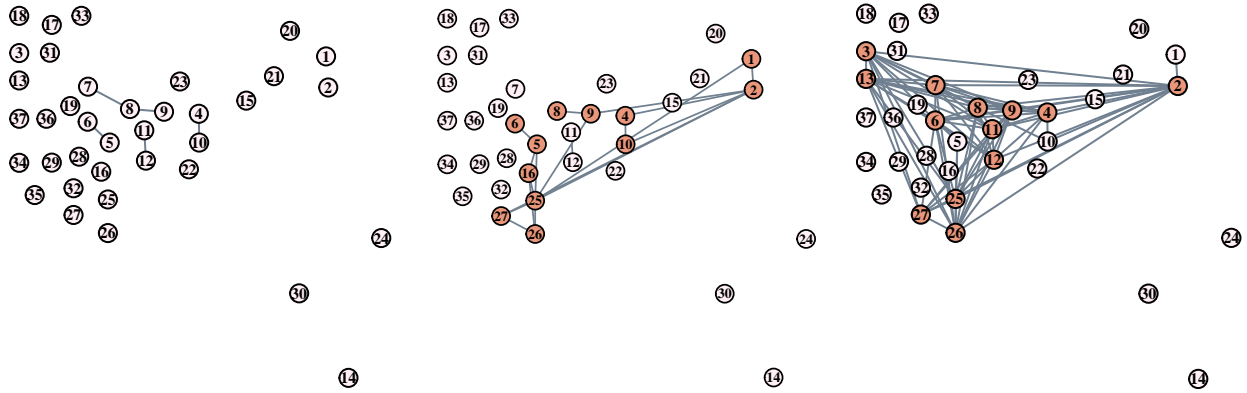


Figure 17: Left: Posterior mode for the representative of cluster 1, with posterior mass 100%, for the last 100,000 iterations. Middle: Posterior mode for the representative of cluster 2, with posterior mass 37%, for the last 100,000 iterations. Right: Posterior mode for the representative of cluster 3, with posterior mass 15%, for the last 100,000 iterations.

among displays due to WiFi connection issues. This remark is also justified by the movements observed in the representative network that correspond to movements among displays within the same building, suggesting that when a user is in a building, the WiFi connection is preserved, and so the application can record the movements of the user.

For cluster $c = 2$, we observe that the posterior mode of the representative network concentrates a relatively smaller posterior mass of 37%, while being slightly denser compared to the representative of cluster 1. Moreover, an SBM structure is revealed, with the mostly connected nodes belonging in the same block. This representative reveals a specific movement pattern of the users allocated in cluster 2, corresponding to the displays located at the central part of the campus (nodes 5, 6, 16, 25, 26, 27), as well as displays located at Infolab (nodes 1 and 2), and Furness College (4 and 10). However, there is a smaller proportion of networks allocated to cluster 2 compared to cluster 1. However, the small false positive (posterior mean 0.02) and high false negative probability (posterior mean 0.49) again indicates that there might be movements of users not recorded by the application.

Lastly, for cluster $c = 3$ the posterior mode of the network representative concentrates a smaller posterior mass compared to the other two representatives equal to 15%. We notice similarities both in the block structure and the connectivity patterns of the representatives of clusters 2 and 3. However, the representative of cluster 3 is notably denser, and some new movements of individuals at displays labelled 3 and 13 (Faraday College and County College) are discovered. We also note that for cluster 3, the posterior means obtained for the false positive and false negative probabilities are similar to cluster 2. In addition, clusters 2 and 3 have a similar proportion of individuals. Overall, a common movement pattern is discovered for individuals in clusters 2 and 3, indicating movements among displays located in the central part of the campus.

We additionally investigate whether users in each cluster interact differently with the displays. Specifically, Figure 18 presents the proportion of times each type of content was shown per cluster of individuals. We have excluded the content "Welcome screen" as it is the default screen introducing the users to the application. We see that individuals assigned to cluster 3 are more active in terms of the range of content shown on the displays visited. This is reasonable considering that the representative of this cluster (Figure 17 right) is the densest, potentially indicating these users are most engaged with the application. For individuals in cluster 1, we see the most common type of content shown is the bus timetable. This can be explained by the movement patterns of the individuals in this cluster 1, as represented by the network in Figure 17 left, which are primarily between displays near to where the bus station is located. Lastly, the type of content mostly shown for individuals of cluster 2 is the weather, which seems to be relevant given the representative movement patterns for these individuals (Figure 17 middle), encompassing longer distances across campus, away from the central covered part.

The analysis of the data assuming $C = 3$ clusters gives sensible and interpretable results as discussed above. The clear allocation of individuals to clusters, and the movement patterns revealed by the cluster specific network representatives, along with the type of content characterising each cluster, indicate that

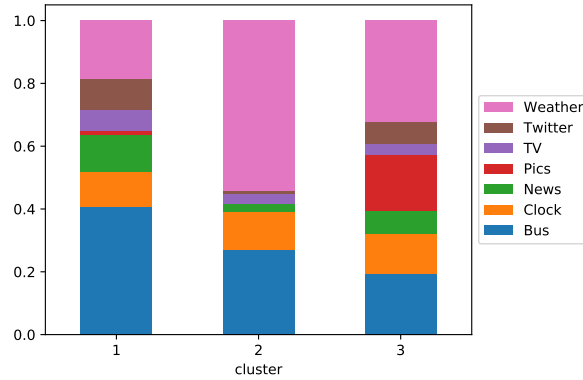


Figure 18: Proportion of times (y axis) each type of content shown (colors) per cluster of individuals (x axis).

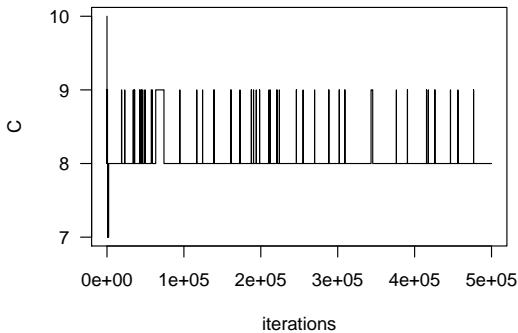


Figure 19: Traceplot for number of clusters detected in each iterations of the MCMC for SFM extension of our model.

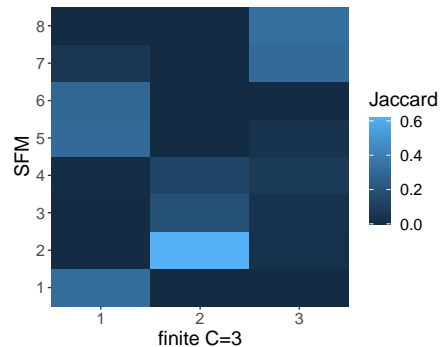


Figure 20: Jaccard similarity between the set of observations in each of the 8 clusters recovered with the SFM extension of our model (y axis) and the set of observations in each of the 3 clusters from our finite mixture model.

three clusters is a sensible choice for the data. However, to investigate the number of clusters further, we implement the SFM extension of our model.

For the SFM, we specify an upper bound of clusters $C_{max} = 10$ and a $\text{Gamma}(1, 400)$ hyperprior to impose a high degree of sparsity. We run the MCMC for 500,000 iterations with a burn-in of 150,000 iterations. In Figure 19, we present the traceplot for the number of clusters in each iteration of the MCMC. To summarise the results from the posterior draws obtained for the cluster membership of the networks, we use the R package **GreedyEPL** (Rastelli and Friel, 2018) that gives a final, optimal partition of the observations. The final partition obtained from the posterior draws of the SFM model suggests the presence of 8 clusters in the data, as can also be seen from the traceplot in Figure 19. This is despite the high sparsity imposed through the hyperprior. However, there are several clusters containing only few observations (between 5 to 14 observations) and only three clusters contain more than 20 observations each, which is not appealing for certain inferences. In particular, one of the main innovations of our model is the ability to infer a cluster specific network representative; however, inferring a network representative for clusters that contain only few observations is not necessarily meaningful.

To compare the final partitions obtained under the finite mixture model with $C = 3$ and the SFM extension, we use a set similarity metric, specifically the Jaccard similarity, which is defined as the size of the intersection divided by the size of the union of the sets, ranging from 0 (no similarity) to 1 (perfect similarity). In Figure 20 we visualise the pairwise comparisons of the partitions obtained under our finite mixture model with $C = 3$ clusters and the SFM model extension of our model using the Jaccard similarity

metric. We notice that the network observations allocated to the three clusters of the finite model, are spread out among the clusters of the SFM extension, with only cluster 2 of the finite model being more similar to cluster 2 of the SFM model.

To account for the cluster configurations obtained across all the iterations of the MCMC, we visualise the proportion of times that each individual is allocated to one of the identified clusters for the SFM, similarly to the visualisation obtained for the finite mixture model in Figure 16 (right). Under the SFM extension of our model, Figure 21 shows a greater variability in the allocation of networks to clusters, without a particularly dominant cluster of networks, compared to the finite mixture model in Figure 16 (right).

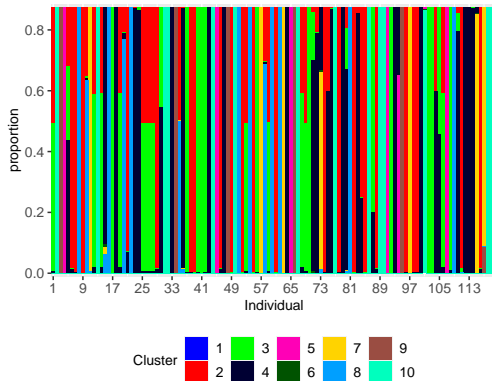


Figure 21: Proportion of times (y axis) an individual’s network (x axis) is allocated to each cluster inferred by the SFM model, after a burn-in of 150,000 iterations.

As noted in Section 5.3, the inference for the number of clusters with the SFM model is highly influenced by the hyperprior of e_0 specified. This also holds for infinite Dirichlet process mixture models, where the posterior distribution of the number of clusters is highly influenced by the precision parameter α of the Dirichlet process, as stated in Frühwirth-Schnatter and Malsiner-Walli (2019). Thus, different hyperprior settings could potentially lead to different conclusions about the number of clusters in the data, with some clusters being hard to interpret in the context of our application, which is a fundamental objective of our modelling framework. We adopt the view that the number of clusters should in part be informed by the applied questions of interest.

Lastly, we formulate our model to allow inferences for matrices of false positive P and false negative Q probabilities with an SBM structure as proposed in Le et al. (2018) and discussed in Section 4.1. The issue arising with this assumption for the Tacita data is the absence of a block structure for the network representative of one of the mixture components as the results showed in this section. Our model makes simultaneous inferences for the block membership of the nodes together with the other model parameters within the MCMC. In contrast, Le et al. (2018) have a two-stage algorithm to infer the block membership of the nodes under the assumption of a single true network in the population. In our setting, this leads to identifiability issues when making inferences for the block specific false positive and false negative probabilities, as suggested by the results of our MCMC, not presented herein due to length limitations. In light of this, exploring alternative MCMC formulations for making inferences for P, Q matrices presents an interesting direction for future work.

6.2 Connectivity patterns in the brain

As described in Section 2.2, this example involves a population of 300 undirected networks corresponding to 10 brain-scans taken for 30 healthy individuals via diffusion magnetic resonance imaging (dMRI). The 300 networks are treated as independent observations, similarly to the studies of Lunagómez et al. (2021) and Arroyo et al. (2021). The nodes of the networks correspond to regions of the brain, and edges denote connections recorded among these regions. Specifically, the network data consist of 200 nodes according to the CC200 atlas (Craddock et al., 2012). Our goal is to identify a cluster of individuals with brain connectivity patterns that differ compared to the majority of the networks in the population, and characterise these with respect to a model parameterisation. We thus implement our outlier cluster detection algorithm, as discussed in Section 4.4.

We initialise our algorithm similarly to the initialisation performed for the Tacita application in Section 6.1. Hence, we determine an initial membership of networks in two different clusters by implementing a k-means algorithm using the R package `kmed` Budiaji (2019). This is done for three distance matrices that correspond to the Jaccard, the wavelets and the l_2 distance metrics. We combine the results using majority vote, to obtain the initial cluster membership of each network. We now consider three distance metrics versus the four distance metrics considered for the Tacita application in Section 6.1, as the prespecified number of clusters is now $C = 2$, and we wish to avoid complications caused by ties. We initialise the network representative by generating its edges through independent Bernoulli draws, with probabilities equal to the proportion of times the corresponding edge is observed in the network data. We initialise the block membership of the representative’s nodes by SBM estimation on the initial representative using the R package `blockmodels` INRA and Leger (2015). For the rest of the parameters of the model we consider three different random initialisations, and run the MCMC for each initialisation for 1,000,000 iterations.

In Figure 22 we present the trace plots of the false positive and false negative probabilities for the majority cluster for all iterations under three different initialisations. In the Supplementary Material, Section 3.2, Figure 11 (Mantziou et al., 2023), we also present the trace plots of the false positive and false negative probabilities for the outlier cluster. We observe that under the three different initialisations the algorithm converges very quickly to the same region. This is encouraging and suggests that a high posterior region has been identified.

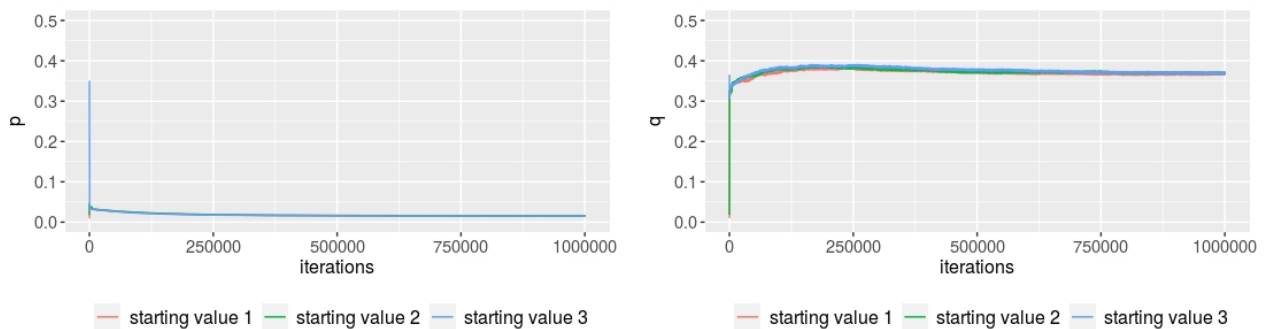


Figure 22: Left: Trace plot for false positive probability p for majority cluster for 1,000,000 iterations and three different initialisations. Right: Trace plot for false negative probability q for majority cluster for 1,000,000 iterations and three different initialisations.

We also compare the results from the three initialisations by obtaining the posterior mode of the representative network for the last 50,000 iterations. In Figure 23 (left) we obtain the posterior mode for the network representative under the first initialisation (posterior mass of 0.08). The colours of the nodes correspond to the block membership. In Figure 23 (middle), we present the not in common edges between the posterior modes of the first and second initialisation to facilitate comparisons. The black edges correspond to the edges present in the posterior mode of the first initialisation and not present in the posterior mode of the second initialisation, and the light gray edges correspond to the edges present in the posterior mode under the second initialisation and not present in the posterior mode of the first initialisation. In Figure 23 (right), we similarly present the not in common edges between the posterior modes of the first and third initialisation. The posterior mode of the second initialisation has posterior mass of 0.08, while the posterior mode of the third initialisation has posterior mass of 0.16.

There are three interesting findings with respect to the representative inferred under the three different initialisations. First, there is only a small proportion of edges not in common among the posterior modes of the three different initialisations, considering the density of the graphs. Second, our algorithm infers the same block structure for the three posterior modes of the representative networks. Third, the posterior masses for the posterior mode representatives are small, but is expected due to the high dimensional space of the networks. In general the results are very encouraging given the size of the networks considered.

We further observe a smaller posterior mean for the false negative probability in the outlier cluster (0.32) compared to the majority cluster (0.37), suggesting that the edges not observed in the network data of the outlier cluster are more likely to be correct compared to the majority cluster. This finding also suggests that the network data in the outlier cluster are sparser compared to the majority cluster. Also, the small posterior means of the false positive probability for both clusters (0.016 for the majority cluster and 0.025

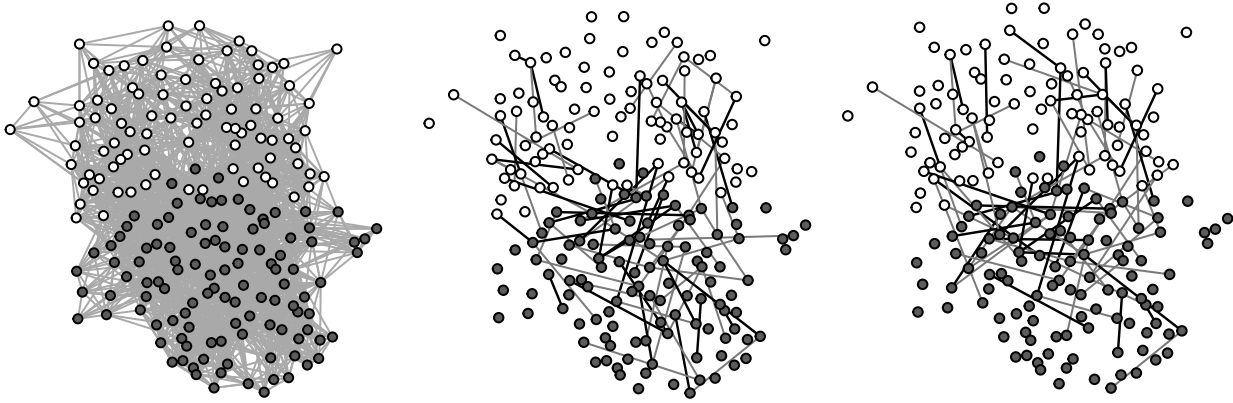


Figure 23: Left: Posterior mode of representative network from 1st initialisation. Middle: Network with not in common edges between posterior modes of representatives from 1st and 2nd initialisation. Right: Network with not in common edges between posterior modes of representatives from 1st and 3rd initialisation. The nodes' colours correspond to the block structure identified under each initialisation.

for the outlier cluster) indicate that the edges observed in the network data are likely correct.

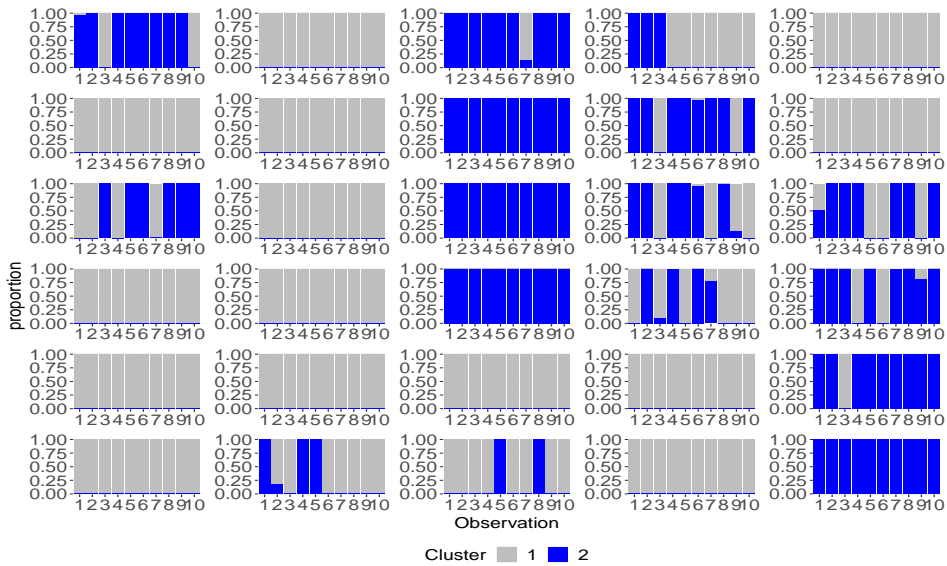


Figure 24: Proportion of times that each of the 10 brain scans from the 30 individuals is allocated in cluster labelled 1 and/or 2.

In Figure 24, we present the results for the cluster membership z of the network observations for the first initialisation. Specifically, we calculate the proportion of times that a network observation is allocated to the majority or the outlier cluster, labelled by 1 or 2 respectively, from the last 100,000 iterations. Each subfigure in Figure 24 shows the cluster allocation of the 10 brain scans obtained for the same individual. We see that our model mostly allocates scans of the same individual to the same cluster. Thus, our model detects similarities among the brain scans of the same individual, giving credence to our model clustering the networks sensibly.

This is a common finding with Arroyo et al. (2021) who performed semi-supervised classification on the brain network population. However, our approach is different to Arroyo et al. (2021) in two ways. First, we implement an unsupervised method to infer underlying clusters of networks. We only pre-define the number of cluster in the population. Second, the interpretation of the results of our model-based clustering method compared to Arroyo et al. (2021) differs significantly, as it reveals a cluster of individuals whose brain connectivity patterns are different to a majority group, and are interpreted through a parametric model.

A common characteristic of human brain networks discussed in the Neuroscience literature, is the exhibition of small-world structures (Bassett and Bullmore, 2006). Networks exhibiting small-worldness are characterised by two main network properties, the clustering coefficient, indicating the level of clique formation in the network, and the average shortest path length, indicating the average length of the shortest paths connecting the nodes. It has been identified that small-worldness is associated with individual cognitive performance (Liao et al., 2017). For example, it has been found that higher intelligence corresponds to shorter path lengths in the brain network (Liao et al., 2017).

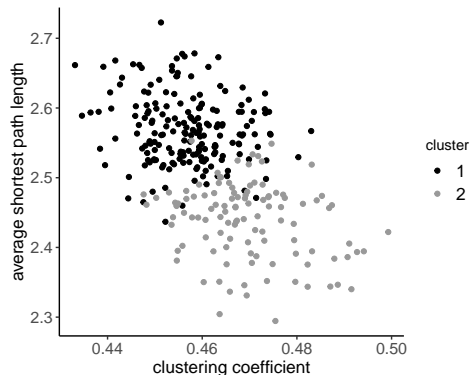


Figure 25: Cluster configuration of brain networks with respect to small-world properties, with x axis corresponding to clustering coefficient and y axis corresponding to average shortest path length. Colours correspond to cluster membership of the brain networks inferred by our algorithm.

In light of this, we investigate how our model clusters the individuals with respect to these two network properties. For each individual brain scan, we calculate the clustering coefficient and the average shortest path length, and use a scatterplot to plot these against each other in Figure 25. The shading of the dots correspond to the cluster membership of each brain scan inferred by our model. We observe a clear distinction between the individual brain scans in each cluster with respect to these two network properties, which are both of interest to the Neuroscience community.

We note here that assuming $C = 2$ clusters, i.e. assuming a single outlier cluster, is appropriate given the objective of our analysis, which is to infer an outlier cluster under a single representative network. Thus, we do not implement the SFM extension of our model in this case. In addition, as the results in this section indicate, the assumption of $C = 2$ gives meaningful and interpretable results, utilising information from the Neuroscience literature.

7 Discussion

In this paper we introduced a mixture model for populations of network data that allows us to identify clusters of networks in a population. To achieve this, we formulated a mixture of measurement error models, and developed a Bayesian framework that allows us to make inferences for all model parameters jointly. This framework permits a diverse specification of the network model for the representative networks in the population, determined according to the type of information we want to exploit based on the data.

Through extensive simulations, we observed our method reliably inferred the model parameters and cluster membership for moderate-sized networks, even for regimes with high noise levels. This is an interesting result, as for high noise levels there is great variability in the structure of the simulated network population, making inference a challenging task. The results suggested that our model can perform well for a range of real data applications. Simulations also examined the model performance for large network and population sizes. This had not been explored by Signorelli and Wit (2020) who also develop a model-based approach for clustering populations of network data. We observed that the absolute errors of the posterior means for the parameters were small, even for large networks and a relatively small sample size. This suggests that our model does not require a large number of observations to make accurate inferences. We also compared our model with the two model-based methods for heterogeneous network populations proposed by Durante et al. (2017) and Signorelli and Wit (2020) respectively. Our model performed similarly to Durante et al. (2017), and better than Signorelli and Wit (2020). Importantly, the parameterisation of our model allows greater interpretability of the results compared to both aforementioned approaches.

We also present an extension to our model that incorporates uncertainty in the number of clusters using the approach suggested in [Frühwirth-Schnatter and Malsiner-Walli \(2019\)](#). The implementation of this approach on simulated data shows that it can accurately recover the true number of clusters as well as the true cluster membership of the simulated networks.

The clustering performed on the Tacita application revealed three different movement patterns of the users. This can be deduced by the network representatives inferred for the clusters, which are primarily characterised by their density. The cluster described by the sparser representative, reveals movements among closely located displays on campus. The second and third cluster identified, enclose denser networks in the population, and the representative of each cluster reveals a specific movement pattern of the individuals therein. As the majority of the networks are very sparse it was encouraging to see that our model was able to separate out the denser networks and further distinguish two different clusters among this subset of individuals. In addition, the type of content shown on the displays visited by the users of each cluster has a meaningful interpretation considering the movement patterns indicated by the network representative in each cluster.

Analysis of the brain network data led to some interesting findings. First, the results suggest we identified a high posterior region for the false positive and false negative probabilities of each cluster. Second, a similar posterior mode for the network representative has been inferred under three different initialisations. These are especially encouraging given the high dimensional space spanned by 200-node network representatives. [Lunagómez et al. \(2021\)](#) who also obtain a network representative in terms of a Fréchet mean for the same brain network population resorted to divide and conquer methods to be able to make inferences, which was not required here. Another interesting finding is that our model partitions the brain networks with respect to two network properties of particular interest to the Neuroscience community.

Our model could potentially incorporate covariates at the node or edge level to inform the inference. In some network applications, it is common to have additional information about the nodes or the edges of the network. For example, the Tacita mobile application also records the type of content shown by the display at the time visited by the user. The incorporation of this additional information could potentially lead to interesting additional findings. It would be also interesting to perform a follow up analysis on data recorded by the Tacita application after the emergence of the COVID-19 pandemic. This would allow us to investigate whether there is a change in the movement patterns of users before and after the pandemic.

In our analysis, the brain network data have been studied as independent network observations, similar to the studies of [Lunagómez et al. \(2021\)](#) and [Arroyo et al. \(2021\)](#). However, in this data set, multiple brain measurements of the same individual are included, and thus the assumption of independence might not be satisfied. The impact of the independence assumption for populations of network observations when dependencies between observations exist has not been studied in the network literature, due to a lack of existing methods for modelling multiple dependent network data. Thus, an interesting direction for future work would be to modify our model in order to capture dependence between network observations.

An interesting result from both the Tacita and the brain networks application, is the high false negative probabilities inferred for the clusters, attributed to the networks' sparsity. Network sparsity is a common issue in many real world network applications. One way to deal with the sparsity would be to consider shrinkage priors e.g. formulating the Horseshoe priors ([Carvalho et al., 2009](#)). Another way to account for network sparsity is to assume the networks are partially observed. In the literature, partially observed networks have been considered under two different perspectives. The coarsening approach focuses on incorporating the coarsening mechanism, that allows us to only partially observe the networks, in the model, and efficiently impute the partially observed data thereafter ([Heitjan and Rubin, 1991](#); [Handcock and Gile, 2010](#); [Heitjan and Rubin, 1990](#); [Kim and Hong, 2012](#)). Another approach focuses on the missingness of certain edges and performs edge prediction ([Koskinen et al., 2013](#); [Marchette and Hohman, 2015](#); [Zhao et al., 2017](#); [Airoldi and Blocker, 2013](#)). Under the first approach, one way we could incorporate the coarsening mechanism is through the assumption of a sampling design, while under the second approach we could have a two-stage method which would first involve performing link prediction, and second performing inference. All the aforementioned approaches require significant modifications of our model and present interesting avenues for future research.

A key challenge arising with the analysis of network data, and especially with data sets consisting of multiple network observations, is the development of methods with good scaling properties as the number of networks' nodes and sample size increase. In our study, we were able to explore the performance of our model for networks with up to 100 nodes and sample sizes of 315 network observations, as well as apply our model and get meaningful results for a Neuroscience application involving networks with 200 nodes and a sample size of 300 observations. The network and sample sizes considered in our study are larger than those

commonly considered in the network literature, and specifically in the analysis of [Durante et al. \(2017\)](#) and [Signorelli and Wit \(2020\)](#). However, there is increasing availability of massive networks for which MCMC approaches like ours are not necessarily practical to implement with respect to the computational cost. The inference of network representatives through MCMC approaches for populations of networks with very large node sets can become very challenging. In such cases, other methods, such as Variational Bayes, could constitute alternatives for making inferences

An additional challenge not considered here is when a population of networks contain both small clusters, together with some much larger clusters. In general, classification techniques are known to struggle when class imbalance is large (i.e. when one class dominates the sample), and so similar challenges could occur in this setting as a result. In particular, correctly inferring cluster membership of small clusters, as well as associated properties of the cluster representatives, could be compromised when the vast majority of networks belong to a single cluster. Addressing this area is an important direction for future research. However, it is encouraging to see that in the brain networks application, small networks clusters are able to be inferred, which indicates the potential of our model to deal with a modest amount of imbalance in cluster sizes.

Our paper contributes to the growing literature on clustering complex types of data. Examples of such studies are the study of [Lu and Marron \(2014\)](#) who focus on identifying clusters of juggling cycles using Functional Principal Component Analysis (FPCA) on diverse data objects, [Song et al. \(2007\)](#) who perform model-based clustering on time-dependent gene expression data using Functional Data Analysis, and [Shen et al. \(2013\)](#) who propose an agglomerative clustering approach for clustering shape data according to their structure. In summary, the flexibility of our modelling framework has been shown to address diverse applied research questions with the potential to be widely applicable to many fields.

Acknowledgment

The authors would like to thank Mateusz Mikusz and Petteri Nurmi for sharing the Tacita data set and for the constructive discussions about the data.

Supplement

Supplement to "Bayesian model-based clustering for populations of network data"

Supplement contains additional details and results for the model, the simulations and the data applications.

Code for Bayesian model-based clustering for populations of network data

This file contains main code for MCMC algorithms, as well as code for simulation experiments and data analysis.

References

- Airoldi, E. M. and Blocker, A. W. (2013). Estimating latent processes on a network from indirect measurements. *Journal of the American Statistical Association*, 108(501):149–164.
- Arroyo, J., Athreya, A., Cape, J., Chen, G., Priebe, C. E., and Vogelstein, J. T. (2021). Inference for multiple heterogeneous networks with a common invariant subspace. *Journal of Machine Learning Research*, 22(142):1–49.
- Arroyo Reli3n, J. D., Kessler, D., Levina, E., and Taylor, S. F. (2019). Network classification with applications to brain connectomics. *The annals of applied statistics*, 13(3):1648.
- Balachandran, P., Kolaczyk, E. D., and Viles, W. D. (2017). On the propagation of low-rate measurement error to subgraph counts in large networks. *The Journal of Machine Learning Research*, 18(1):2025–2057.
- Bassett, D. S. and Bullmore, E. (2006). Small-world brain networks. *The neuroscientist*, 12(6):512–523.
- Budiaji, W. (2019). *kmed: Distance-Based K-Medoids*. R package version 0.3.0.
- Carvalho, C. M., Polson, N. G., and Scott, J. G. (2009). Handling sparsity via the horseshoe. In *Artificial Intelligence and Statistics*, pages 73–80. PMLR.

- Chang, J., Kolaczyk, E. D., and Yao, Q. (2020). Estimation of subgraph densities in noisy networks. *Journal of the American Statistical Association*, pages 1–14.
- Chatterjee, S. (2015). Matrix estimation by universal singular value thresholding. *The Annals of Statistics*, 43(1):177–214.
- Craddock, R. C., James, G. A., Holtzheimer III, P. E., Hu, X. P., and Mayberg, H. S. (2012). A whole brain fmri atlas generated via spatially constrained spectral clustering. *Human brain mapping*, 33(8):1914–1928.
- Diquigiovanni, J. and Scarpa, B. (2019). Analysis of association football playing styles: An innovative method to cluster networks. *Statistical Modelling*, 19(1):28–54.
- Donnat, C. and Holmes, S. (2018). Tracking network dynamics: A survey using graph distances. *The Annals of Applied Statistics*, 12(2):971–1012.
- Durante, D., Dunson, D. B., and Vogelstein, J. T. (2017). Nonparametric bayes modeling of populations of networks. *Journal of the American Statistical Association*.
- Fields, S. and Song, O.-k. (1989). A novel genetic system to detect protein–protein interactions. *Nature*, 340(6230):245–246.
- Frühwirth-Schnatter, S. and Malsiner-Walli, G. (2019). From here to infinity: sparse finite versus dirichlet process mixtures in model-based clustering. *Advances in data analysis and classification*, 13(1):33–64.
- Ginestet, C. E., Li, J., Balachandran, P., Rosenberg, S., and Kolaczyk, E. D. (2017). Hypothesis testing for network data in functional neuroimaging. *The Annals of Applied Statistics*, pages 725–750.
- Gollini, I. and Murphy, T. B. (2016). Joint modeling of multiple network views. *Journal of Computational and Graphical Statistics*, 25(1):246–265.
- Handcock, M. S. and Gile, K. J. (2010). Modeling social networks from sampled data. *The Annals of Applied Statistics*, 4(1):5.
- Heitjan, D. F. and Rubin, D. B. (1990). Inference from coarse data via multiple imputation with application to age heaping. *Journal of the American Statistical Association*, 85(410):304–314.
- Heitjan, D. F. and Rubin, D. B. (1991). Ignorability and coarse data. *The annals of statistics*, pages 2244–2253.
- Hoff, P. D., Raftery, A. E., and Handcock, M. S. (2002). Latent space approaches to social network analysis. *Journal of the american Statistical association*, 97(460):1090–1098.
- INRA and Leger, J.-B. (2015). *blockmodels: Latent and Stochastic Block Model Estimation by a 'V-EM' Algorithm*. R package version 1.1.1.
- Jiang, X., Gold, D., and Kolaczyk, E. D. (2011). Network-based auto-probit modeling for protein function prediction. *Biometrics*, 67(3):958–966.
- Josephs, N., Li, W., and Kolaczyk, E. D. (2021). Network recovery from unlabeled noisy samples. In *2021 55th Asilomar Conference on Signals, Systems, and Computers*, pages 1268–1273. IEEE.
- Keribin, C. (2000). Consistent estimation of the order of mixture models. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 49–66.
- Kim, H. and Park, H. (2007). Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23(12):1495–1502.
- Kim, J. K. and Hong, M. (2012). Imputation for statistical inference with coarse data. *Canadian Journal of Statistics*, 40(3):604–618.
- Kolaczyk, E. D., Lin, L., Rosenberg, S., Walters, J., and Xu, J. (2020). Averages of unlabeled networks: Geometric characterization and asymptotic behavior. *The Annals of Statistics*, 48(1):514–538.

- Koskinen, J. H., Robins, G. L., Wang, P., and Pattison, P. E. (2013). Bayesian analysis for partially observed network data, missing ties, attributes and actors. *Social Networks*, 35(4):514–527.
- Le, C. M., Levin, K., and Levina, E. (2018). Estimating a network from multiple noisy realizations. *Electronic Journal of Statistics*, 12(2):4697–4740.
- Le, C. M. and Li, T. (2020). Linear regression and its inference on noisy network-linked data. *arXiv preprint arXiv:2007.00803*.
- Levin, K., Athreya, A., Tang, M., Lyzinski, V., Park, Y., and Priebe, C. E. (2017). A central limit theorem for an omnibus embedding of multiple random graphs and implications for multiscale network inference. *arXiv preprint arXiv:1705.09355*.
- Li, W., Sussman, D. L., and Kolaczyk, E. D. (2021). Causal inference under network interference with noise. *arXiv preprint arXiv:2105.04518*.
- Liao, X., Vasilakos, A. V., and He, Y. (2017). Small-world human brain networks: perspectives and challenges. *Neuroscience & Biobehavioral Reviews*, 77:286–300.
- Lu, X. and Marron, J. (2014). Analysis of juggling data: Object oriented data analysis of clustering in acceleration functions. *Electronic Journal of Statistics*, 8(2):1842–1847.
- Lunagómez, S., Olhede, S. C., and Wolfe, P. J. (2021). Modeling network populations via graph distances. *Journal of the American Statistical Association*, 116(536):2023–2040.
- Malsiner-Walli, G., Frühwirth-Schnatter, S., and Grün, B. (2016). Model-based clustering based on sparse finite gaussian mixtures. *Statistics and computing*, 26(1):303–324.
- Mantziou, A., Lunagómez, S., and Mitra, R. (2023). Supplement to "bayesian model-based clustering for populations of network data". *DOI*.
- Marchette, D. J. and Hohman, E. L. (2015). Utilizing covariates in partially observed networks. In *2015 18th International Conference on Information Fusion (Fusion)*, pages 166–172. IEEE.
- Mukherjee, S. S., Sarkar, P., and Lin, L. (2017). On clustering network-valued data. *Advances in neural information processing systems*.
- Neal, R. M. (2000). Markov chain sampling methods for dirichlet process mixture models. *Journal of computational and graphical statistics*, 9(2):249–265.
- Newman, M. E. (2018). Estimating network structure from unreliable measurements. *Physical Review E*, 98(6):062321.
- Nielsen, A. M. and Witten, D. (2018). The multiple random dot product graph model. *arXiv preprint arXiv:1811.12172*.
- Peixoto, T. P. (2018). Reconstructing networks with unknown and heterogeneous errors. *Physical Review X*, 8(4):041011.
- Prasad, G., Joshi, S. H., Nir, T. M., Toga, A. W., Thompson, P. M., (ADNI, A. D. N. I., et al. (2015). Brain connectivity and novel network measures for alzheimer’s disease classification. *Neurobiology of aging*, 36:S121–S131.
- Priebe, C. E., Sussman, D. L., Tang, M., and Vogelstein, J. T. (2015). Statistical inference on errorfully observed graphs. *Journal of Computational and Graphical Statistics*, 24(4):930–953.
- Rastelli, R. and Friel, N. (2018). Optimal bayesian estimators for latent variable cluster models. *Statistics and Computing*, 28(6):1169–1186.
- Richardson, S. and Green, P. J. (1997). On bayesian analysis of mixtures with an unknown number of components (with discussion). *Journal of the Royal Statistical Society: series B (statistical methodology)*, 59(4):731–792.

- Richiardi, J., Eryilmaz, H., Schwartz, S., Vuilleumier, P., and Van De Ville, D. (2011). Decoding brain states from fmri connectivity graphs. *Neuroimage*, 56(2):616–626.
- Schütze, H., Manning, C. D., and Raghavan, P. (2008). *Introduction to information retrieval*, volume 39. Cambridge University Press Cambridge.
- Shaw, P., Mikusz, M., Nurmi, P., and Davies, N. (2018). Tacita: A privacy preserving public display personalisation service. In *Proceedings of the 2018 ACM International Joint Conference and 2018 International Symposium on Pervasive and Ubiquitous Computing and Wearable Computers*, pages 448–451.
- Shen, W., Wang, Y., Bai, X., Wang, H., and Latecki, L. J. (2013). Shape clustering: Common structure discovery. *Pattern Recognition*, 46(2):539–550.
- Signorelli, M. and Wit, E. C. (2020). Model-based clustering for populations of networks. *Statistical Modelling*, 20(1):9–29.
- Song, J. J., Lee, H.-J., Morris, J. S., and Kang, S. (2007). Clustering of time-course gene expression data using functional data analysis. *Computational biology and chemistry*, 31(4):265–274.
- Wang, S., Arroyo, J., Vogelstein, J. T., and Priebe, C. E. (2019). Joint embedding of graphs. *IEEE transactions on pattern analysis and machine intelligence*.
- White, J. G., Southgate, E., Thomson, J. N., and Brenner, S. (1986). The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Philos Trans R Soc Lond B Biol Sci*, 314(1165):1–340.
- Young, J.-G., Cantwell, G. T., and Newman, M. (2020). Bayesian inference of network structure from unreliable data. *Journal of Complex Networks*, 8(6):cnaa046.
- Young, S. J. and Scheinerman, E. R. (2007). Random dot product graph models for social networks. pages 138–149.
- Zhang, J., Cheng, W., Wang, Z., Zhang, Z., Lu, W., Lu, G., and Feng, J. (2012). Pattern classification of large-scale functional brain networks: identification of informative neuroimaging markers for epilepsy. *PloS one*, 7(5):e36733.
- Zhao, Y., Wu, Y.-J., Levina, E., and Zhu, J. (2017). Link prediction for partially observed networks. *Journal of Computational and Graphical Statistics*, 26(3):725–733.
- Zuo, X.-N., Anderson, J. S., Bellec, P., Birn, R. M., Biswal, B. B., Blautzik, J., Breitner, J. C., Buckner, R. L., Calhoun, V. D., Castellanos, F. X., et al. (2014). An open science resource for establishing reliability and reproducibility in functional connectomics. *Scientific data*, 1(1):1–13.

Supplement to "Bayesian model-based clustering for populations of network data"

Anastasia Mantziou¹, Simón Lunagómez², and Robin Mitra³

¹The Alan Turing Institute

²Departamento de Estadística, ITAM

³Department of Statistical Science, University College London

In this document we provide supplementary material to the article "Bayesian model-based clustering for multiple network data". In Section 1 we provide more details about the MCMC scheme introduced in Section 4 of the main article. In Section 2 we provide additional results from the simulation studies performed in Section 5 of the main article. In Section 3, we provide additional results for the real data applications discussed in Section 6 of the main article. In Section 4 we provide details of the MCMC algorithm implemented.

1 Additional details for the MCMC scheme

In this Section we provide the full conditional posteriors for the model parameters that are updated through a Gibbs sampler, as discussed in Section 4.3 of the main article.

The full conditional posterior for the probability of a network to belong to cluster τ is given by

$$P(\tau | \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z}, \mathbf{w}, \mathbf{b}, \boldsymbol{\theta}, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) \propto P(\mathbf{z} | \tau) \cdot P(\tau | \boldsymbol{\psi})$$

where $P(\mathbf{z} | \tau) = \text{Multinomial}(1; \tau_1, \dots, \tau_C)$ and $P(\tau | \boldsymbol{\psi}) = \text{Dirichlet}(\boldsymbol{\psi})$, as specified in Section 4.2 of the main article. Thus, we have

$$\begin{aligned} P(\tau | \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z}, \mathbf{w}, \mathbf{b}, \boldsymbol{\theta}, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) &\propto \prod_{j=1}^N \tau_{z_j} \cdot \Gamma(\psi C) \cdot \Gamma(\boldsymbol{\psi})^{-C} \cdot \prod_{c=1}^C \tau_c^{\psi-1} \\ &\propto \Gamma(\psi C) \cdot \Gamma(\boldsymbol{\psi})^{-C} \cdot \tau_{z_1} \cdots \tau_{z_N} \cdot \tau_1^{\psi-1} \cdots \tau_C^{\psi-1} \\ &\propto \Gamma(\psi C) \cdot \Gamma(\boldsymbol{\psi})^{-C} \cdot \tau_1^{\eta_1} \cdots \tau_C^{\eta_C} \cdot \tau_1^{\psi-1} \cdots \tau_C^{\psi-1} \propto \tau_1^{\eta_1+\psi-1} \cdots \tau_C^{\eta_C+\psi-1} \end{aligned} \quad (1)$$

where $\eta_c = \sum_{j=1}^N 1_c(z_j)$, $c = 1, \dots, C$, denotes the number of network data that belong to cluster c . Thence we obtain,

$$P(\tau | \mathbf{A}_{\mathcal{G}^*}, \mathbf{p}, \mathbf{q}, \mathbf{z}, \mathbf{w}, \mathbf{b}, \boldsymbol{\theta}, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) = \text{Dirichlet}(\boldsymbol{\psi} + \boldsymbol{\eta}_1, \dots, \boldsymbol{\psi} + \boldsymbol{\eta}_C).$$

The derivation of the full conditional posterior for the vector of the nodes' block-membership probabilities \mathbf{w}_c for cluster c , is similar to the derivation of the the full conditional posterior for τ , as already described above, thus we have

$$P(\mathbf{w}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \boldsymbol{\theta}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) \propto P(\mathbf{b}_c | \mathbf{w}_c) \cdot P(\mathbf{w}_c | \boldsymbol{\chi})$$

where $P(\mathbf{b}_c | \mathbf{w}_c) = \text{Multinomial}(\mathbf{w}_c)$ and $P(\mathbf{w}_c | \boldsymbol{\chi}) = \text{Dirichlet}(\boldsymbol{\chi})$, as specified in Section 4.1 of the main article. Hence we obtain

$$\begin{aligned} P(\mathbf{w}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \boldsymbol{\theta}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) &\propto \prod_{i=1}^n w_{c,b_i} \cdot \Gamma(\chi K) \cdot \Gamma(\boldsymbol{\chi})^{-K} \cdot \prod_{k=1}^K w_{c,k}^{\chi-1} \\ &\propto \Gamma(\chi K) \cdot \Gamma(\boldsymbol{\chi})^{-K} \cdot w_{c,b_1} \cdots w_{c,b_n} \cdot w_{c,1}^{\chi-1} \cdots w_{c,K}^{\chi-1} \\ &\propto \Gamma(\chi K) \cdot \Gamma(\boldsymbol{\chi})^{-K} \cdot w_{c,1}^{h_1} \cdots w_{c,K}^{h_K} \cdot w_{c,1}^{\chi-1} \cdots w_{c,K}^{\chi-1} \propto w_{c,1}^{(h_1+\chi)-1} \cdots w_{c,K}^{(h_K+\chi)-1}. \end{aligned}$$

where h_k denotes the number of the nodes that belong to block k . Thus the full conditional posterior for \mathbf{w}_c is

$$P(\mathbf{w}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \boldsymbol{\theta}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) = \text{Dirichlet}(\chi + h_1, \dots, \chi + h_K).$$

The full conditional posterior for the vector of the block-specific probabilities of an edge occurrence, $\boldsymbol{\theta}_c$, for the network representative of cluster c is

$$P(\boldsymbol{\theta}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \mathbf{w}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) \propto P(A_{\mathcal{G}_c^*} | \mathbf{w}_c, \mathbf{b}_c, \boldsymbol{\theta}_c) \cdot P(\boldsymbol{\theta}_c | \boldsymbol{\epsilon}, \boldsymbol{\zeta})$$

where $P(A_{\mathcal{G}_c^*} | \mathbf{w}_c, \mathbf{b}_c, \boldsymbol{\theta}_c) = \text{SBM}(\mathbf{w}_c, \mathbf{b}_c, \boldsymbol{\theta}_c)$ and $P(\boldsymbol{\theta}_c | \boldsymbol{\epsilon}, \boldsymbol{\zeta}) = \text{Beta}(\boldsymbol{\epsilon}, \boldsymbol{\zeta})$, as specified in Section 4.1 of the main article. Thus,

$$\begin{aligned} & P(\boldsymbol{\theta}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \mathbf{w}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) \\ & \propto \prod_{(i,j):i<j} \theta_{c,b_i b_j}^{A_{\mathcal{G}_c^*}(i,j)} (1 - \theta_{c,b_e, i b_e, j})^{1 - A_{\mathcal{G}_c^*}(i,j)} \cdot \prod_{k=1}^K \prod_{l=1}^K \theta_{c,kl}^{\epsilon_{kl}-1} (1 - \theta_{c,kl})^{\zeta_{kl}-1} \\ & \propto \prod_{k=1}^K \prod_{l=1}^K \theta_{c,kl}^{A_{\mathcal{G}_c^*}[kl]} (1 - \theta_{c,kl})^{n_{c,kl} - A_{\mathcal{G}_c^*}[kl]} \theta_{c,kl}^{\epsilon_{kl}-1} (1 - \theta_{c,kl})^{\zeta_{kl}-1} \\ & \propto \prod_{k=1}^K \prod_{l=1}^K \theta_{c,kl}^{A_{\mathcal{G}_c^*}[kl] + \epsilon_{kl} - 1} (1 - \theta_{c,kl})^{n_{c,kl} - A_{\mathcal{G}_c^*}[kl] + \zeta_{kl} - 1}, \end{aligned}$$

where $A_{\mathcal{G}_c^*}[kl] = \sum_{(i,j):b_{c,i}=k, b_{c,j}=l} A_{\mathcal{G}_c^*}(i,j)$ represents the sum of the entries for the pairs of nodes of the network representative for cluster c that have block membership k, l respectively, and $n_{c,kl} = \sum_{(i,j):i \neq j} \mathbb{1}(b_{c,i} = k, b_{c,j} = l)$ is the number of the pairs of nodes of the representative of cluster c that have membership k, l accordingly. Hence we obtain

$$P(\boldsymbol{\theta}_c | A_{\mathcal{G}_c^*}, p_c, q_c, \mathbf{z}, \boldsymbol{\tau}, \mathbf{b}_c, \mathbf{w}_c, A_{\mathcal{G}_1}, \dots, A_{\mathcal{G}_N}) = \text{Beta}(A_{\mathcal{G}_c^*}[kl] + \epsilon_{kl}, \zeta_{kl} + n_{c,kl} - A_{\mathcal{G}_c^*}[kl]).$$

2 Additional details for the Simulation Studies

2.1 Additional details for simulation study for moderate-sized networks

In this Section we provide the results for the simulation regimes presented in Section 5.1 of the main article (Table 1 in the main article). Specifically, in Tables 1-13 we present the posterior means and credible intervals for the false positive probabilities p_c , false negative probabilities q_c , and block specific edge probabilities $\boldsymbol{\theta}_c$. In Tables 14-15 we present the posterior means for the probability of a node to belong to a block \mathbf{w}_c .

In addition, Tables 16-17 show the proportion of times that the Hamming distance between the true representatives and the posterior representatives is less than or equal to 1, 5 and 10 respectively, for each simulation regime. Table 18 shows the mean clustering entropy and mean clustering purity calculated for each simulation regime, as discussed in Section 5.1 of the main article.

		SBM ₁		SBM ₂	
p_c	q_c	posterior mean of p_c	posterior mean of q_c	posterior mean of p_c	posterior mean of q_c
0.1	0.2	(0.09,0.10,0.11)	(0.20,0.21,0.21)	(0.10,0.10,0.10)	(0.21,0.19,0.19)
	0.3	(0.10,0.10,0.10)	(0.31,0.29,0.30)	(0.11,0.10,0.10)	(0.30,0.31,0.30)
0.2	0.1	(0.19,0.19,0.20)	(0.10,0.10,0.11)	(0.20,0.20,0.20)	(0.10,0.09,0.10)
	0.3	(0.20,0.20,0.19)	(0.29,0.29,0.30)	(0.20,0.20,0.19)	(0.30,0.32,0.30)
0.3	0.1	(0.31,0.30,0.31)	(0.10,0.10,0.10)	(0.31,0.30,0.31)	(0.10,0.10,0.10)
	0.2	(0.30,0.29,0.30)	(0.20,0.19,0.21)	(0.30,0.30,0.28)	(0.20,0.19,0.20)

Table 1: Posterior means for false positive probabilities p_c and false negative probabilities q_c , for $c \in \{1, 2, 3\}$.

SBM ₁				
p_c	q_c	credible interval for p_1	credible interval for p_2	credible interval for p_3
0.1	0.2	(0.09,0.10)	(0.09,0.11)	(0.10,0.11)
	0.3	(0.09,0.10)	(0.10,0.11)	(0.09,0.11)
0.2	0.1	(0.18,0.20)	(0.18,0.20)	(0.19,0.21)
	0.3	(0.19,0.21)	(0.19,0.21)	(0.18,0.20)
0.3	0.1	(0.29,0.32)	(0.29,0.31)	(0.30,0.32)
	0.2	(0.29,0.32)	(0.28,0.31)	(0.29,0.31)

Table 2: 95 % credible intervals for false positive probabilities p_c , for $c \in \{1, 2, 3\}$, under SBM 1.

SBM ₂				
p_c	q_c	credible interval for p_1	credible interval for p_2	credible interval for p_3
0.1	0.2	(0.09,0.11)	(0.09,0.11)	(0.09,0.11)
	0.3	(0.10,0.11)	(0.09,0.10)	(0.09,0.11)
0.2	0.1	(0.19,0.21)	(0.19,0.21)	(0.19,0.21)
	0.3	(0.19,0.21)	(0.19,0.21)	(0.18,0.20)
0.3	0.1	(0.30,0.32)	(0.29,0.31)	(0.30,0.32)
	0.2	(0.29,0.31)	(0.29,0.31)	(0.27,0.30)

Table 3: 95 % credible intervals for false positive probabilities p_c , for $c \in \{1, 2, 3\}$, under SBM 2.

SBM ₁				
p_c	q_c	credible interval for q_1	credible interval for q_2	credible interval for q_3
0.1	0.2	(0.19,0.21)	(0.20,0.21)	(0.19,0.22)
	0.3	(0.30,0.32)	(0.28,0.30)	(0.28,0.31)
0.2	0.1	(0.09,0.10)	(0.09,0.11)	(0.10,0.12)
	0.3	(0.28,0.30)	(0.28,0.31)	(0.28,0.31)
0.3	0.1	(0.09,0.11)	(0.09,0.11)	(0.10,0.11)
	0.2	(0.19,0.21)	(0.18,0.20)	(0.20,0.22)

Table 4: 95 % credible intervals for false negative probabilities q_c , for $c \in \{1, 2, 3\}$, under SBM 1.

SBM ₂				
p_c	q_c	credible interval for q_1	credible interval for q_2	credible interval for q_3
0.1	0.2	(0.20,0.22)	(0.18,0.20)	(0.18,0.20)
	0.3	(0.29,0.31)	(0.30,0.32)	(0.29,0.32)
0.2	0.1	(0.09,0.11)	(0.09,0.10)	(0.10,0.11)
	0.3	(0.29,0.31)	(0.30,0.33)	(0.29,0.31)
0.3	0.1	(0.09,0.11)	(0.09,0.11)	(0.09,0.11)
	0.2	(0.19,0.21)	(0.18,0.20)	(0.19,0.21)

Table 5: 95 % credible intervals for false negative probabilities q_c , for $c \in \{1, 2, 3\}$, under SBM 2.

SBM ₁				
p_c	q_c	posterior means of $(\theta_{11}^{(1)}, \theta_{12}^{(1)}, \theta_{22}^{(1)})$	posterior means of $(\theta_{11}^{(2)}, \theta_{12}^{(2)}, \theta_{22}^{(2)})$	posterior means of $(\theta_{11}^{(3)}, \theta_{12}^{(3)}, \theta_{22}^{(3)})$
0.1	0.2	(0.84,0.22,0.90)	(0.83,0.21,0.82)	(0.79,0.18,0.74)
	0.3	(0.84,0.22,0.90)	(0.83,0.21,0.82)	(0.79,0.18,0.74)
0.2	0.1	(0.84,0.22,0.90)	(0.83,0.21,0.82)	(0.79,0.18,0.74)
	0.3	(0.84,0.22,0.90)	(0.83,0.21,0.82)	(0.74,0.18,0.79)
0.3	0.1	(0.84,0.22,0.90)	(0.82,0.21,0.83)	(0.74,0.18,0.79)
	0.2	(0.90,0.22,0.84)	(0.82,0.21,0.83)	(0.79,0.18,0.74)

Table 6: Posterior means for block specific edge probabilities θ_c , for $c \in \{1, 2, 3\}$, under SBM 1.

SBM ₂				
p_c	q_c	posterior means of $(\theta_{11}^{(1)}, \theta_{12}^{(1)}, \theta_{22}^{(1)})$	posterior means of $(\theta_{11}^{(2)}, \theta_{12}^{(2)}, \theta_{22}^{(2)})$	posterior means of $(\theta_{11}^{(3)}, \theta_{12}^{(3)}, \theta_{22}^{(3)})$
0.1	0.2	(0.68,0.02,0.78)	(0.60,0.05,0.71)	(0.78,0.03,0.83)
	0.3	(0.68,0.02,0.78)	(0.60,0.05,0.71)	(0.78,0.03,0.83)
0.2	0.1	(0.68,0.02,0.78)	(0.60,0.05,0.71)	(0.78,0.03,0.83)
	0.3	(0.68,0.02,0.78)	(0.71,0.05,0.60)	(0.83,0.03,0.78)
0.3	0.1	(0.68,0.02,0.78)	(0.60,0.05,0.71)	(0.78,0.03,0.83)
	0.2	(0.68,0.02,0.78)	(0.71,0.05,0.60)	(0.83,0.03,0.78)

Table 7: Posterior means for block specific edge probabilities θ_c , for $c \in \{1, 2, 3\}$, under SBM 2.

SBM ₁				
p_c	q_c	credible interval for $\theta_{11}^{(1)}$	credible interval for $\theta_{12}^{(1)}$	credible interval for $\theta_{22}^{(1)}$
0.1	0.2	(0.73,0.94)	(0.15,0.30)	(0.82,0.97)
	0.3	(0.73,0.94)	(0.15,0.30)	(0.82,0.97)
0.2	0.1	(0.73,0.94)	(0.15,0.30)	(0.82,0.97)
	0.3	(0.73,0.94)	(0.15,0.30)	(0.82,0.97)
0.3	0.1	(0.73,0.94)	(0.15,0.30)	(0.82,0.97)
	0.2	(0.82,0.97)	(0.15,0.30)	(0.73,0.94)

Table 8: 95 % credible intervals for block specific edge probabilities θ_1 of cluster 1, under SBM 1.

SBM ₂				
p_c	q_c	credible interval for $\theta_{11}^{(1)}$	credible interval for $\theta_{12}^{(1)}$	credible interval for $\theta_{22}^{(1)}$
0.1	0.2	(0.59,0.77)	(0.00,0.04)	(0.58,0.96)
	0.3	(0.60,0.77)	(0.00,0.04)	(0.59,0.96)
0.2	0.1	(0.60,0.77)	(0.00,0.04)	(0.59,0.96)
	0.3	(0.60,0.77)	(0.00,0.04)	(0.59,0.96)
0.3	0.1	(0.60,0.77)	(0.00,0.04)	(0.59,0.96)
	0.2	(0.59,0.77)	(0.00,0.04)	(0.58,0.96)

Table 9: 95 % credible intervals for block specific edge probabilities θ_1 of cluster 1, under SBM 2.

SBM ₁				
p_c	q_c	credible interval for $\theta_{11}^{(2)}$	credible interval for $\theta_{12}^{(2)}$	credible interval for $\theta_{22}^{(2)}$
0.1	0.2	(0.73,0.92)	(0.14,0.29)	(0.70,0.92)
	0.3	(0.73,0.92)	(0.14,0.29)	(0.70,0.92)
0.2	0.1	(0.73,0.92)	(0.14,0.29)	(0.70,0.92)
	0.3	(0.73,0.92)	(0.14,0.29)	(0.70,0.92)
0.3	0.1	(0.70,0.92)	(0.14,0.29)	(0.73,0.92)
	0.2	(0.70,0.92)	(0.14,0.29)	(0.73,0.92)

Table 10: 95 % credible intervals for block specific edge probabilities θ_2 of cluster 2, under SBM 1.

SBM ₂				
p_c	q_c	credible interval for $\theta_{11}^{(2)}$	credible interval for $\theta_{12}^{(2)}$	credible interval for $\theta_{22}^{(2)}$
0.1	0.2	(0.47,0.72)	(0.01,0.09)	(0.58,0.83)
	0.3	(0.47,0.72)	(0.01,0.09)	(0.58,0.83)
0.2	0.1	(0.47,0.72)	(0.01,0.09)	(0.58,0.83)
	0.3	(0.47,0.72)	(0.01,0.09)	(0.58,0.83)
0.3	0.1	(0.47,0.73)	(0.01,0.09)	(0.57,0.83)
	0.2	(0.47,0.72)	(0.01,0.09)	(0.58,0.83)

Table 11: 95 % credible intervals for block specific edge probabilities θ_2 of cluster 2, under SBM 2.

SBM ₁				
p_c	q_c	credible interval for $\theta_{11}^{(3)}$	credible interval for $\theta_{12}^{(3)}$	credible interval for $\theta_{22}^{(3)}$
0.1	0.2	(0.68,0.90)	(0.11,0.26)	(0.63,0.85)
	0.3	(0.68,0.90)	(0.11,0.26)	(0.63,0.85)
0.2	0.1	(0.68,0.90)	(0.12,0.26)	(0.63,0.85)
	0.3	(0.68,0.90)	(0.12,0.26)	(0.63,0.85)
0.3	0.1	(0.68,0.90)	(0.12,0.26)	(0.63,0.85)
	0.2	(0.68,0.90)	(0.11,0.26)	(0.63,0.85)

Table 12: 95 % credible intervals for block specific edge probabilities θ_3 of cluster 3, under SBM 1.

SBM ₂				
p_c	q_c	credible interval for $\theta_{11}^{(3)}$	credible interval for $\theta_{12}^{(3)}$	credible interval for $\theta_{22}^{(3)}$
0.1	0.2	(0.59,0.96)	(0.00,0.06)	(0.75,0.89)
	0.3	(0.58,0.96)	(0.00,0.06)	(0.75,0.89)
0.2	0.1	(0.59,0.96)	(0.00,0.06)	(0.75,0.89)
	0.3	(0.59,0.96)	(0.00,0.06)	(0.75,0.89)
0.3	0.1	(0.59,0.96)	(0.00,0.06)	(0.75,0.90)
	0.2	(0.59,0.96)	(0.00,0.06)	(0.75,0.89)

Table 13: 95 % credible intervals for block specific edge probabilities θ_3 of cluster 3, under SBM 2.

SBM ₁				
p_c	q_c	posterior means of $(w_1^{(1)}, w_2^{(1)})$	posterior means of $(w_1^{(2)}, w_2^{(2)})$	posterior means of $(w_1^{(3)}, w_2^{(3)})$
0.1	0.2	(0.48,0.52)	(0.52,0.48)	(0.48,0.52)
	0.3	(0.48,0.52)	(0.52,0.48)	(0.48,0.52)
0.2	0.1	(0.48,0.52)	(0.52,0.48)	(0.48,0.52)
	0.3	(0.48,0.52)	(0.52,0.48)	(0.52,0.48)
0.3	0.1	(0.48,0.52)	(0.48,0.52)	(0.52,0.48)
	0.2	(0.52,0.48)	(0.48,0.52)	(0.48,0.52)

Table 14: Posterior means for the probability of a node to belong to a block w_c for $c \in \{1, 2, 3\}$, under SBM 1.

SBM ₂				
p_c	q_c	posterior means of $(w_1^{(1)}, w_2^{(1)})$	posterior means of $(w_1^{(2)}, w_2^{(2)})$	posterior means of $(w_1^{(3)}, w_2^{(3)})$
0.1	0.2	(0.70,0.30)	(0.52,0.48)	(0.30,0.70)
	0.3	(0.70,0.30)	(0.52,0.48)	(0.30,0.70)
0.2	0.1	(0.70,0.30)	(0.52,0.48)	(0.30,0.70)
	0.3	(0.70,0.30)	(0.48,0.52)	(0.30,0.70)
0.3	0.1	(0.70,0.30)	(0.52,0.48)	(0.30,0.70)
	0.2	(0.70,0.30)	(0.48,0.52)	(0.30,0.70)

Table 15: Posterior means for the probability of a node to belong to a block w_c for $c \in \{1, 2, 3\}$, under SBM 2.

SBM ₁				
p_c	q_c	$d_H \leq 1$	$d_H \leq 5$	$d_H \leq 10$
0.1	0.2	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
	0.3	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
0.2	0.1	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
	0.3	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
0.3	0.1	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
	0.2	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)

Table 16: Proportion of times that the Hamming distance between the posterior representatives and the true representatives is less or equal than 1, 5 and 10 respectively, under SBM 1.

SBM ₂				
p_c	q_c	$d_H \leq 1$	$d_H \leq 5$	$d_H \leq 10$
0.1	0.2	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
	0.3	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
0.2	0.1	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
	0.3	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
0.3	0.1	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)
	0.2	(1.00,1.00,1.00)	(1.00,1.00,1.00)	(1.00,1.00,1.00)

Table 17: Proportion of times that the Hamming distance between the posterior representatives and the true representatives is less or equal than 1, 5 and 10 respectively, under SBM 2.

SBM ₁		SBM ₂			
p_c	q_c	Mean Entropy	Mean Purity	Mean Entropy	Mean Purity
0.1	0.2	0	1	0	1
	0.3	0	1	0	1
0.2	0.1	0	1	0	1
	0.3	0	1	0	1
0.3	0.1	0	1	0	1
	0.2	0	1	0	1

Table 18: Mean clustering entropy and clustering purity

2.2 Additional details for simulations involving varying network sizes and population sizes

In Figures 1-4, we illustrate the representatives of clusters 2 and 3, with 25, 50, 75 and 100 nodes respectively, generated for the simulation study presented in Section 5.2 of the main article.

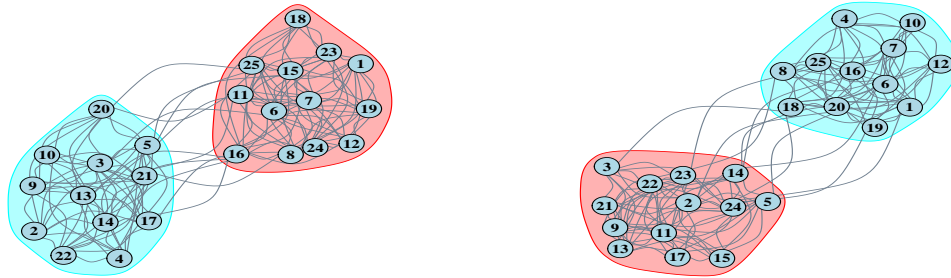


Figure 1: 25-node representatives of cluster 2 (left) and cluster 3 (right).

Similarly to the posterior summaries obtained with moderate network sizes, we obtain the proportion of times that the Hamming distance between posterior draws of the representative and the true representative is less than or equal to 1, 5 and 10, for each cluster. We consider the final 350,000 posterior draws after a burn-in of 150,000 iterations. The results are summarised in Figure 5. The multiple subfigures correspond to the 25, 50, 75 and 100 node representatives.

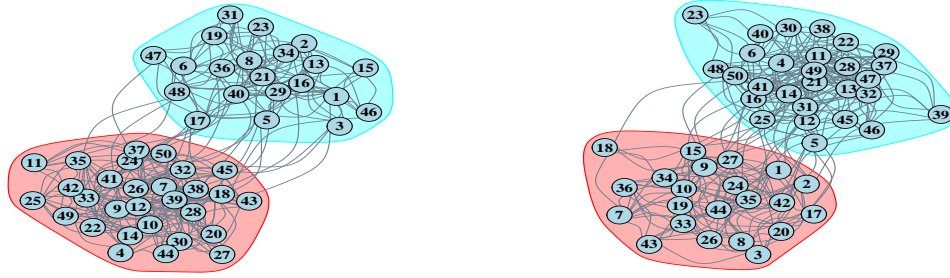


Figure 2: 50-node representatives of cluster 2 (left) and cluster 3 (right).

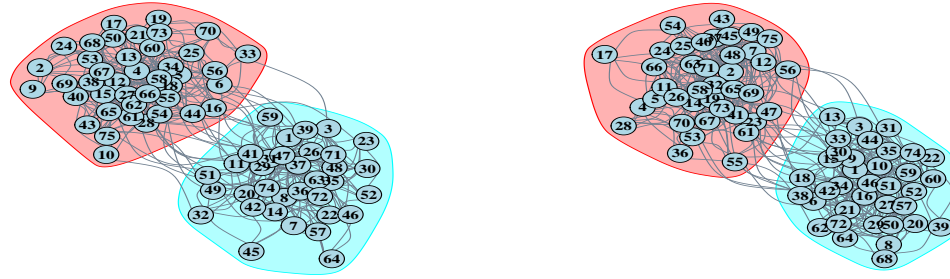


Figure 3: 75-node representatives of cluster 2 (left) and cluster 3 (right).

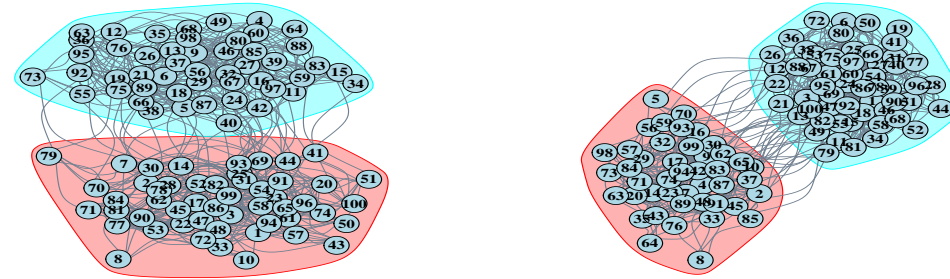


Figure 4: 100-node representatives of cluster 2 (left) and cluster 3 (right).

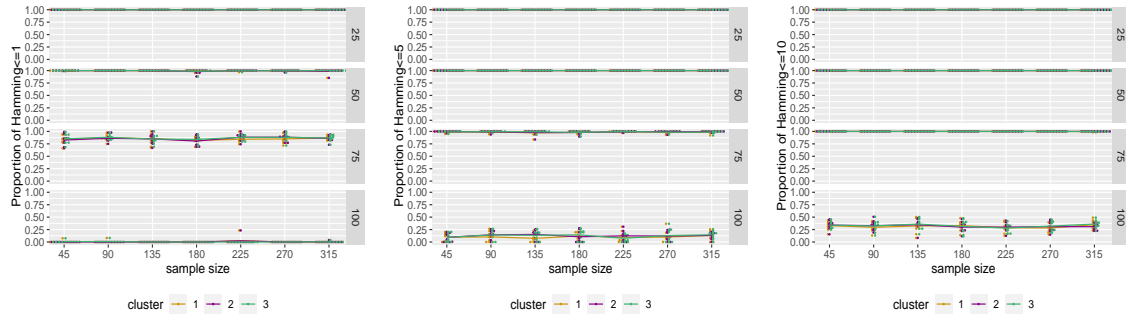


Figure 5: Left: Proportion of times the Hamming distance is less or equal to 1 (y axis) for 25-node, 50-node, 75-node and 100-node network representatives and varying population sizes (x axis). Middle: Proportion of times the Hamming distance is less or equal to 5 (y axis) for 25-node, 50-node, 75-node and 100-node network representatives and varying population sizes (x axis). Right: Proportion of times the Hamming distance is less or equal to 10 (y axis) for 25-node, 50-node, 75-node and 100-node network representatives and varying population sizes (x axis).

3 Additional details for real data application

3.1 Additional details for the analysis in Section 6.1

In this section we present details of the analysis performed on the data collected by the Tacita mobile application monitoring the movement of individuals across a University campus. We present the results of the exploratory data analysis (EDA) conducted, as well as additional results from fitting our mixture model to the data. Table 19 presents the labels of the nodes corresponding to each display located on Lancaster University campus.

Display name	Node label	Display name	Node label
SCC (C-floor)	1	ISS	20
Infolab Foyer	2	Pendle College	21
Faraday Left	3	New Engineering	22
Engineering Foyer (far)	4	Fylde College	23
LZ1	5	Graduate College	24
LZ3	6	Library A	25
Furness 1	7	Library B	26
Furness 2	8	Library C	27
Furness College	9	Bowland Main B	28
Engineering Foyer (near)	10	Bowland North B	29
LEC 1	11	Hotel Conference	30
LEC 2	12	Chemistry A	31
County College	13	Psychology	32
Lonsdale College	14	Physics	33
Grizedale College	15	Law 2	34
The Base	16	Law 1	35
Faraday B	17	Welcome Screen 1	36
Faraday C	18	Welcome Screen 2	37
Bowland JCR	19		

Table 19: Node label assigned to each display on Lancaster University campus.

As discussed in the main article, we performed EDA on the Tacita multiple network data through the use of network distance metrics. Specifically, for each distance metric, we derive a distance matrix that encloses the pairwise distances of the networks in the population. Then, the (i, j) element of a distance matrix corresponds to the distance between graphs \mathcal{G}_i and \mathcal{G}_j , for the specified distance metric. We consider various distance metrics, as different metrics can give us different information with respect to the presence of clusters in the network population. We consider the Hamming, the Jaccard, the l_2 distance and the distance based on wavelets. To graphically represent the distance matrices for each distance metric, we use a Multidimensional Scaling (MDS) plot. The MDS algorithm maps objects in a 2-d space, respecting their pairwise distances. In Figures 6 and 7, we plot the MDS representation under the distance matrices obtained under the Hamming, the Jaccard, the l_2 , and the wavelets distance metrics.

We also present some additional results after fitting our mixture model to the Tacita data. Figures 8 and 9 show the trace plots of 400,000 posterior draws for the false positive probabilities, p_c , and false negative probabilities, q_c , with $c \in \{1, 2, 3\}$, after a burn-in of 100,000 iterations. In addition, Figure 10 shows the proportion of times that each of the 37 nodes of the representative of cluster $c = 1$ is allocated to either block 1 or 2.

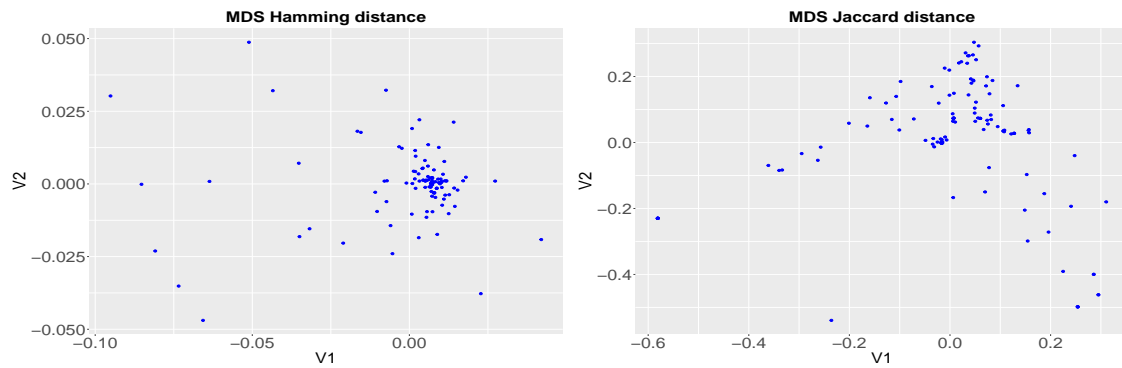


Figure 6: Left: MDS for Hamming distance matrix for the Tacita data. Right: MDS for Jaccard distance matrix for the Tacita data.

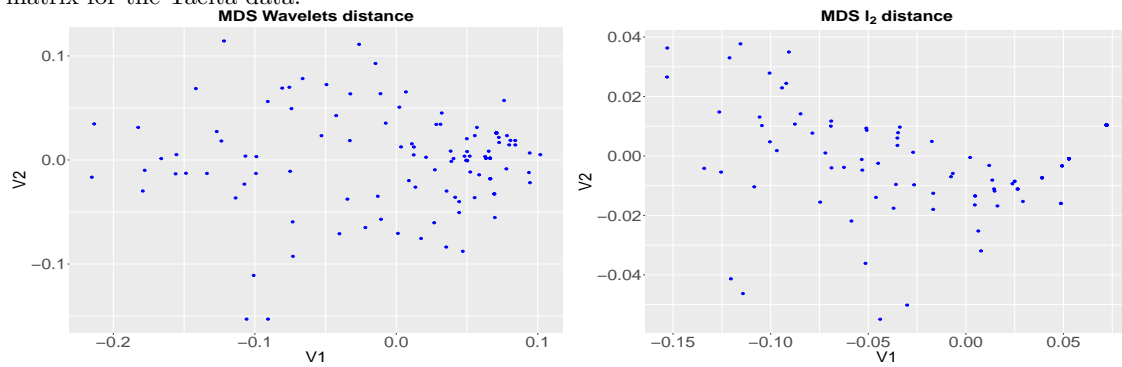


Figure 7: Left: MDS for wavelets distance matrix for the Tacita data. Right: MDS for l_2 distance matrix for the Tacita data.

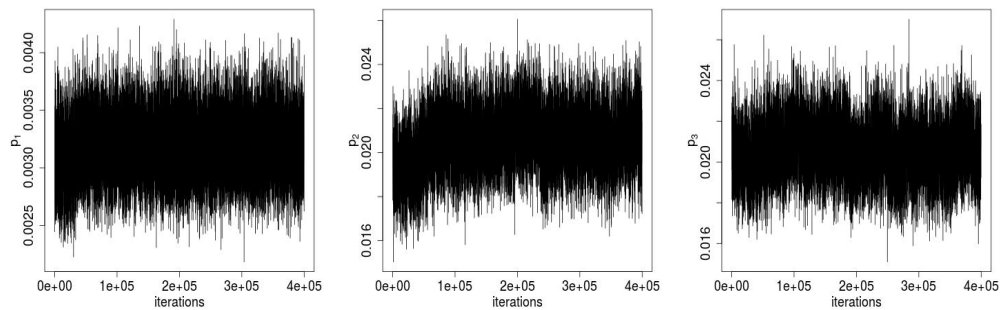


Figure 8: Trace plots of false positive probabilities p_c for $c = \{1, 2, 3\}$, for 400,000 iterations of the MCMC after a burn-in of 100,000 iterations.

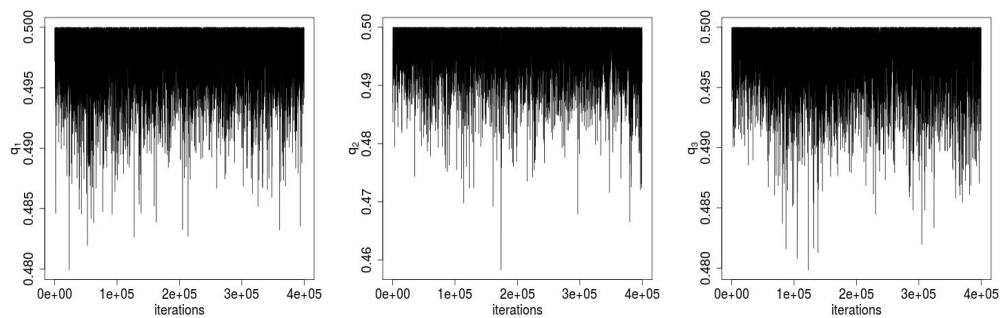


Figure 9: Trace plots of false negative probabilities q_c for $c = \{1, 2, 3\}$, for 400,000 iterations of the MCMC after a burn-in of 100,000 iterations.

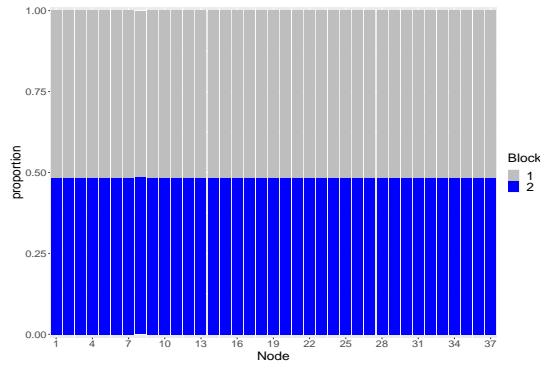


Figure 10: Proportion of times that each node (x axis) of representative of cluster 1, is allocated to each of the two blocks, after a burn-in of 100,000 iterations.

3.2 Additional details for the analysis in Section 6.2

In Figure 11 we show the trace plots for the false positive, p_{out} , and false negative, q_{out} , probabilities, for the outlier cluster of networks detected, under three different initialisations, and 1,000,000 iterations of the MCMC.

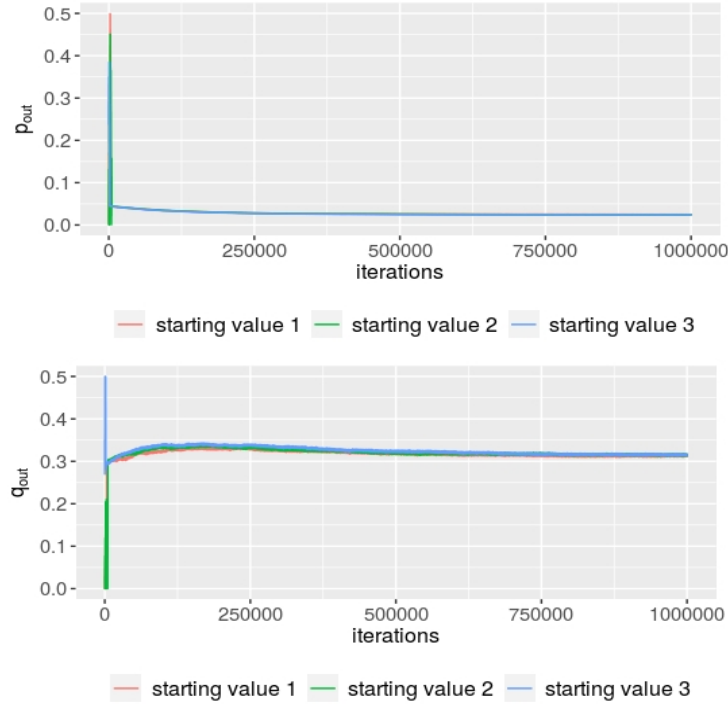


Figure 11: Top: Trace plot of false positive probability for outlier cluster p_{out} for 1,000,000 iterations and three different initialisations. Bottom: Trace plot of false negative probability for outlier cluster q_{out} for 1,000,000 iterations and three different initialisations.

4 MCMC algorithm

Below, we present details of how the MCMC algorithm is implemented to make posterior inferences from the proposed model. In Section 4.5 of the main article, we further discuss the Sparse Finite Mixture extension to the algorithm presented herein.

Algorithm 1: MCMC Algorithm for Clustering Network Populations

Input: $\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_N; C, K, M, w_0, \theta_0, \alpha_0, \beta_0, \gamma_0, \delta_0, \epsilon_0, \zeta_0, \psi, \chi$
Output: Posterior distributions of $A_{\mathcal{G}_1^*}, \dots, A_{\mathcal{G}_C^*}, p_1, \dots, p_C, q_1, \dots, q_C, \tau_1, \dots, \tau_C, z_1, \dots, z_N,$
 $\theta_1, \dots, \theta_C, \mathbf{w}_1, \dots, \mathbf{w}_C, \mathbf{b}_1, \dots, \mathbf{b}_C$
Initialisation: starting values $A_{\mathcal{G}_1^*}^{(0)}, \dots, A_{\mathcal{G}_C^*}^{(0)}, p_1^{(0)}, \dots, p_C^{(0)}, q_1^{(0)}, \dots, q_C^{(0)}, \tau_1^{(0)}, \dots, \tau_C^{(0)},$
 $z_1^{(0)}, \dots, z_N^{(0)}, \theta_1^{(0)}, \dots, \theta_C^{(0)}, \mathbf{w}_1^{(0)}, \dots, \mathbf{w}_C^{(0)}, \mathbf{b}_1^{(0)}, \dots, \mathbf{b}_C^{(0)}$
for $i \leftarrow 1$ **to** M **do**
 Gibbs step: Update τ_1, \dots, τ_C
 compute: $\eta_c = \sum_{j=1}^N \mathbb{1}_c(z_j^{(i-1)})$ for $c = 1, \dots, C$
 sample: $\tau_1^{(i)}, \dots, \tau_C^{(i)} \sim \text{Dir}(\psi + \eta_1, \dots, \psi + \eta_C)$
 for $c \leftarrow 1$ **to** C **do**
 MH step with a mixture of kernels: Update $A_{\mathcal{G}_c^*}$ or p_c or q_c
 sample: $v \sim \text{Multinomial}(\xi_1, \dots, \xi_L)$
 Depending on the value of v , update either $A_{\mathcal{G}_c^*}$ or p_c or q_c as per the Measurement Error model with SBM structure, where the sum in likelihood is over the networks $\{j : z_j^{(i-1)} = c\}$
 Gibbs step: Update \mathbf{w}_c
 compute: $h_k = \sum_{j=1}^n \mathbb{1}_k(b_j^{(i-1)})$
 sample: $\mathbf{w}_c^{(i)} \sim \text{Dir}(\chi + h_1, \dots, \chi + h_K)$
 Gibbs step: Update θ_c
 compute: $A[st] = \sum_{(u,v): b_u=s, b_v=t} A_{\mathcal{G}_c^*}^{(i)}(u, v)$ and
 $n_{st} = \sum_{(u,v): u \neq v} \mathbb{1}(b_u = s, b_v = t)$ for $s, t \in \{1, \dots, K\}$
 sample: $\theta_{c,st}^{(i)} \sim \text{Beta}(A[st] + \epsilon_0, \zeta_0 + n_{st} - A[st])$
 Gibbs step: Update \mathbf{b}_c
 for $j \leftarrow 1$ **to** n **do**
 compute: $p_{kj} = w_{c,k}^{(i)} \cdot \prod_{m=1}^n \theta_{kb_{c,m}^{(i-1)}}^{(i)A(j,m)} (1 - \theta_{kb_{c,m}^{(i-1)}}^{(i)})^{1-A(j,m)}$ for $k = 1, \dots, K$
 sample: $b_{c,j}^{(i)} \sim \text{Multin}(p_{1j}, \dots, p_{Kj})$
 Gibbs step: Update z_1, \dots, z_N
 for $j \leftarrow 1$ **to** N **do**
 compute: $p_{cj} = \tau_c^{(i)} \cdot \prod_{(u,v): u < v} \left((1 - q_c^{(i)})^{A_{\mathcal{G}_j}(u,v)} q_c^{(i)(1-A_{\mathcal{G}_j}(u,v))} \right)^{A_{\mathcal{G}_c^*}^{(i)}(u,v)}$
 $\left(p_c^{(i)A_{\mathcal{G}_j}(u,v)} (1 - p_c^{(i)})^{1-A_{\mathcal{G}_j}(u,v)} \right)^{1-A_{\mathcal{G}_c^*}^{(i)}(u,v)}$ for $c = 1, \dots, C$
 sample: $z_j^{(i)} \sim \text{Multin}(p_{1j}, \dots, p_{Cj})$
