# Mixture of multilayer stochastic block models for multiview clustering

**Kylliann De Santiago**                         KYLLIANN.DESANTIAGO@UNIV-EVRY.FR
*Université Paris-Saclay, CNRS, Univ Evry,*
*Laboratoire de Mathématiques et Modélisation d'Evry,*
*91037, Évry-Courcouronnes, France.*

**Marie Szafranski**                             MARIE.SZAFRANSKI@UNIV-EVRY.FR
*ENSIIE, 91025, Évry-Courcouronnes, France.*
*Université Paris-Saclay, CNRS, Univ Evry,*
*Laboratoire de Mathématiques et Modélisation d'Evry,*
*91037, Évry-Courcouronnes, France.*

**Christophe Ambroise**                          CHRISTOPHE.AMBROISE@UNIV-EVRY.FR
*Université Paris-Saclay, CNRS, Univ Evry,*
*Laboratoire de Mathématiques et Modélisation d'Evry,*
*91037, Évry-Courcouronnes, France.*

## Abstract

In this work, we propose an original method for aggregating multiple clustering coming from different sources of information. Each partition is encoded by a co-membership matrix between observations. Our approach uses a mixture of multilayer Stochastic Block Models (SBM) to group co-membership matrices with similar information into components and to partition observations into different clusters, taking into account their specificities within the components. The identifiability of the model parameters is established and a variational Bayesian EM algorithm is proposed for the estimation of these parameters. The Bayesian framework allows for selecting an optimal number of clusters and components. The proposed approach is compared using synthetic data with consensus clustering and tensor-based algorithms for community detection in large-scale complex networks. Finally, the method is utilized to analyze global food trading networks, leading to structures of interest.

**Keywords:** Stochastic Block Model, Multiview clustering, Multilayer Network, Bayesian Framework, Integrated Classification Likelihood

## 1 Introduction

Most everyday learning situations are achieved by integrating different sources of information, such as vision, touch and hearing. A source of information in a given format will be referred to as a modality or a *view*. Multimodal or multiview machine learning aims to learn models from multiple views (e.g. text, sound, image, etc.) in order to represent, translate, align, fusion, or co-learn (see Zhao et al., 2017; Baltrušaitis et al., 2018; Cornuéjols et al., 2018, for instance).

Graphs provide a powerful and intuitive way to represent complex systems of relationships between individuals. They provide an effective and informative representation of the

system. Constructing graphs from each view allows to use graph machine learning for multimodal clustering (Ektefaie et al., 2023).

In clustering framework, the output of algorithms are often a partition or a membership matrix $\mathbf{Z}$. This information, although useful, has the drawback of strongly depending on the number of clusters chosen when using the algorithm. To avoid this problem, $\mathbf{Z}$ can be transform into an adjacency matrice $\mathbf{A}$ with

$$A_{ij} = \begin{cases} 1, & \text{if individuals } i, j \text{ are linked in the same cluster,} \\ 0, & \text{otherwise.} \end{cases}$$

The process of combining numerous data clusters that have already been discovered using various clustering algorithms or approaches is known as *meta clustering*. Finding connections and similarities across clusters that might not be immediately obvious when looking at them separately is the aim of meta or *consensus clustering* (Monti et al., 2003; Li et al., 2015; Liu et al., 2018). Model-based consensus clustering offer advantages: knowing the redundancy of information sources and their complementarity, obtaining a final clustering from all the outputs already carried out, allowing the best possible grouping of individuals through different information sources. Moreover, model-based consensus clustering allows to have an evaluation criterion on the performance of the model (e.g. log-likelihood, evidence, etc.) and, at least in the Bayesian framework, criteria for model selection (Biernacki et al., 2010).

The corresponding learning models vary based on their fusion strategy. The three main categories of methods are early, intermediate, and late fusion of views. Late fusion is well suited to clustering since each view is often associated to dedicated efficient clustering algorithms.

**Contribution.** In this work we propose to estimate a coordinated representation produced by learning separate clustering for each view and then coordination through a probabilistic model: MIxture of Multilayer Integrator Stochastic Block Model (mimi-SBM). Our model is a Bayesian mixture of multilayer SBM that takes into account several sources of information, and where the membership clustering is traversing as illustrated in Figure 1.

In simpler terms, each individual belongs exclusively to one group, and not to a multitude of groups across views or sources. In the context of meta-clustering, this has the advantage of allowing the model to find common information for each group, by striving for clustering redundancy across sources. Moreover, by applying a mixture model on the views, we can take into account the particularities of each source of information, and define the redundant and complementary information sources in order to draw a maximum of information from them. Finally, with the Bayesian framework, the development of a model selection criterion, both for the mixture of views and the number of clusters is possible by deriving it from the evidence lower bound. The identifiability of the model parameters is established and a variational Bayesian EM algorithm is proposed for the estimation of these parameters.

**Organization of the paper.** The paper is structured as follows: Firstly, we delve into the related works in more detail, providing a selective survey of the existing literature and research in the field. Next, we present the description of our innovative *mimi-SBM* model, outlining its key components and parameter estimation. We then focus on model selection and variational parameter initialization, where we compare different criteria to determine
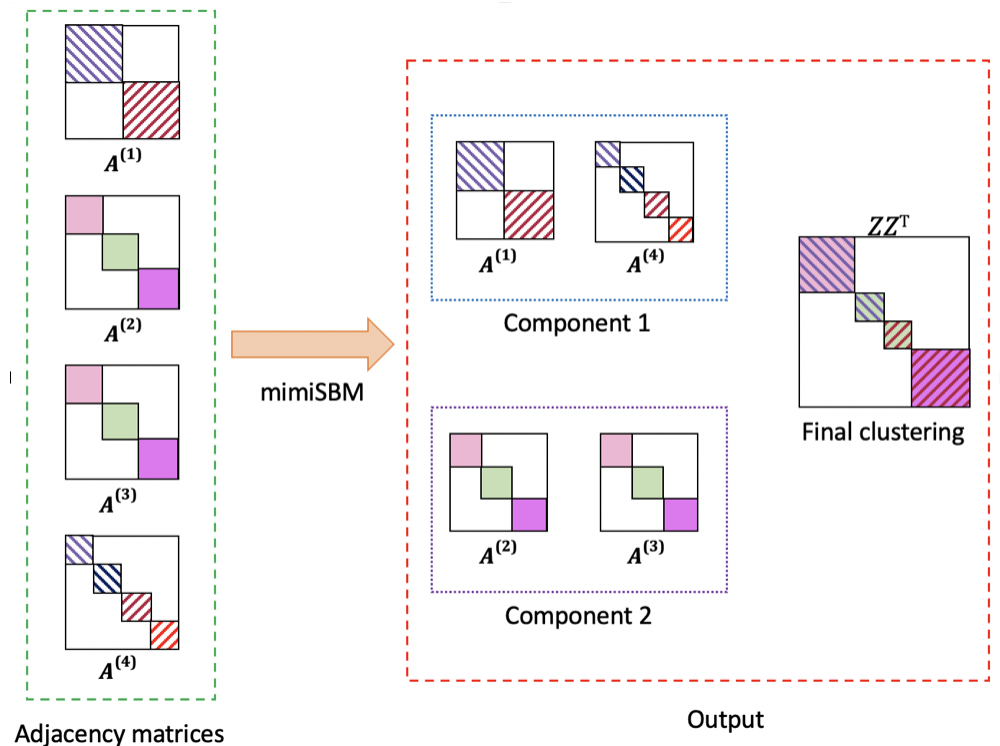
Figure 1: Illustration of mimi-SBM. Left: Four adjacency matrices $\mathbf{A}^{(1)}, \cdots, \mathbf{A}^{(4)}$ coming from four different views organized into two components. Right: identification of the two components from the views (local and complementary information) and clustering of the observations described by the classification matrix $\mathbf{Z}$ (global and consensus information).

the most effective approach. To evaluate the performance of our proposed approach, we conduct synthetic experiments that allow for a thorough comparison and evaluation. Finally, we engage in a discussion about the results obtained, and their implications, and provide insights into potential future directions for further research.

## 2 Background

In multiview clustering, various fusion strategies have been devised to effectively combine information from different sources. One possible taxonomy of these strategies is related to the timing of fusion: early, intermediate, or late. *Early fusion* starts by merging the different views before clustering. *Intermediate fusion* considers the integration of views within the clustering algorithm. *Late fusion* consists of clustering each view separately and then integrating all the resulting partitions into a unique integrated clustering. This paper focuses on late multiview clustering. In this scenario, all layers represented as adjacency matrices collectively form a tensor.

This latter strategy harnesses the benefits of employing specific and well-suited clustering approaches for each view and takes advantage of their complementarity and redundancies by merging their results in a subsequent phase. In this context, *consensus clustering* serves as

a baseline and will be first described in this section. On the other hand, *latent* or *stochastic block models* provide advantageous characteristics for multiview clustering in terms of model selection. In this section, we also provide a focused overview of these approaches, to which our method, tailored to late fusion clustering, belongs.

## 2.1 Consensus clustering

Consensus clustering (Monti et al., 2003; Fred and Jain, 2005), also known as cluster ensemble (Strehl and Ghosh, 2002; Golalipour et al., 2021), is a technique used to find a single partition from multiple clustering solutions.

It is used to integrate and analyze multiple clustering results obtained from different algorithms, parameter settings, or subsets of the data. It aims to find a consensus or agreement among the individual clustering solutions to obtain a more robust and reliable clustering result. Consensus clustering is like asking for multiple opinions (clusters) and then finding a common answer (consensus) that represents an overall agreement.

Each run generates a set of clusters, which can be represented as a partition matrix where each entry indicates the cluster assignment of each data point. An agreement matrix is constructed from all the partition matrices. The entries in this matrix indicate the frequency with which a pair of data points co-occur in the same cluster across all solutions. The final clusters are determined by applying a clustering algorithm on the agreement matrix.

In particular, the Monte Carlo reference-based consensus clustering (M3C) (Monti et al., 2003) combines multiple clustering solutions generated by applying different clustering algorithms and parameters to the same dataset using random (re)sampling techniques.

## 2.2 Block models for multiview clustering

In the framework of block models for multiview clustering, the views are commonly denoted as a collection of $V$ *graphs* or more often as $V$ *layers* within a single network, and the terminologies of *multigraph*, *multilayer* or even *multiplex* may be employed.

With a wide expanse of literature existing on this subject, we narrow our focus on studies in which the distinct views represent varying types of interactions among a common set of $N$ observations. However, it's important to note that our scope excludes studies that aim to establish partitions with overlaps or mixed memberships, as well as those where views show specific dependencies (spatial or temporal for instance). Works of this nature can be discovered in the references provided below.

### 2.2.1 Multilayer SBM

Multilayer SBM (MLSBM) approaches are focused on identifying a partition $\mathbf{Z}$ with $K$ blocks of observations that encompass the different layers. A popular kind of inference for block model estimation relies on Variational Expectation Maximization (VEM) algorithms (Daudin et al., 2008). Besides, when the data originate from a MLSBM, different estimation techniques can also be applied to identify the partition. Out of these, spectral clustering finds widespread usage. (Von Luxburg, 2007; Von Luxburg et al., 2008).

**VEM inference.** In this setting, different SBM based approaches have been proposed for multilayer (Han et al., 2015; Paul and Chen, 2016) or similarly for multiplex (Barbillon et al.,

2017) networks. Han et al. (2015) propose a consistent maximum-likelihood estimate (MLE) and explore the asymptotic properties of class memberships when the number of relations grows. Paul and Chen (2016) also study the consistency of two other MLEs when the number of nodes or the number of types of edges grow together. Barbillon et al. (2017) introduce an Erdös-Rényi model that may also integrate covariates related to pairs of observations and use an Integrated Completed Likelihood (ICL) Biernacki et al. (2010) for the purpose of model selection. Also based on VEM estimation, the work of Boutalbi et al. (2021) is grounded on Latent Block Models (LBM).

**Spectral clustering.** Here, we shed light on several extensive research efforts focused on spectral clustering under the assumption of data generated by a Multilayer Stochastic Block Model. Han et al. (2015) investigate the asymptotic characteristics related to spectral clustering. Chen and Hero (2017) introduce a framework for multilayer spectral graph clustering that includes a scheme for adapting layer weights, while also offering statistical guarantees for the reliability of clustering. Mercado et al. (2018) presents a spectral clustering algorithm for multilayer graphs that relies on the matrix power mean of Laplacians. Paul and Chen (2020) show the consistency of the global optimizers of co-regularized spectral clustering and also for the orthogonal linked matrix factorization. Finally, Huang et al. (2022) propose integrated spectral clustering methods based on convex layer aggregations.

### 2.2.2 Multiway block models

In contrast to the works presented above, where the aim is to establish a partition across the observations, the approaches presented below focus on establishing multiway structures, especially between- and within-layer partitions. In this context, the MLSBM can evolve into either a Mixture of Multilayer SBM (MMLSBM) or expand into a Tensor Block Model (TBM), depending on the specific research communities and their focus.

**Mixture of multilayer SBM.** Stanley et al. (2016) introduced one of the first approaches that integrated a multilayer SBM with a mixture of layers, using a two-step greedy inference method. In the initial step, it infers a SBM for each layer and groups together SBMs with similar parameters. In the second step, these outcomes serve as the starting point for an iterative procedure that simultaneously identifies $Q$ strata spanning the $V$ layers. In each stratum $s$, the nodes are independently distributed into $K_s$ blocks, leading to $Q$ membership matrices $\{\mathbf{Z}^1, \cdots, \mathbf{Z}^Q\}$.

In pursuit of the same goal, Fan et al. (2022) develop an alternating minimization algorithm that offers theoretical guarantees for both between-layer and within-layer clustering errors. Rebafka (2023) proposes a Bayesian framework for a finite mixture of MLSBM and employs a hierarchical agglomerative algorithm for the clustering process. It initiates with individual singleton clusters and then progressively merges clusters of networks according to an ICL criterion also used for model selection.

Also, Pensky and Wang (2021) presents a versatile model for diverse multiplex networks, including both MLSBM and MMLSBM. Note that in the latter scenario, they make the assumption that the number of blocks within each group of layers remains consistent, such that $K_s = K$, $\forall s$. They perform a spectral clustering on the layers and then aggregate the resulting block connectivity matrices to determine the between-layer partition of observations. Using this model as a foundation, Noroozi and Pensky (2022) introduces a more efficient

resolution technique rooted in sparse subspace clustering (Elhamifar and Vidal, 2013). They demonstrate that this algorithm consistently achieves strong between-layer clustering results. In both studies, they provide valuable insights comparing their assumption that $K_s = K$ for all components with the scenario where $K_s$ is considered a known value for each component $s$. This discussion is particularly relevant in the context of methods that are not designed for the task of model selection.

**Tensor block models.** A different strategy for addressing the challenge of late fusion multiview clustering involves tensorial modeling and estimation techniques. Wang and Zeng (2019) position their research within the context of higher-order tensors. They introduce a least-square estimation method for (sparse) TBM and demonstrate the reliability of block structure recovery as the data tensor's dimension increases by providing consistency guarantees. Han et al. (2022) suggest employing high-order spectral clustering as an initialization of a high-order Lloyd algorithm. They establish convergence guarantees and statistical optimality under the assumption of sub-Gaussian noise.

Boutalbi et al. (2020) introduce an extension of Latent Block Models to handle tensors. They consider multivariate normal distributions for continuous data and Bernoulli distributions for categorical data, implementing a VEM algorithm for this purpose.

Finally, Jing et al. (2021) employs the Tucker decomposition to conduct alternating regularized low-rank approximations of the tensor. This technique consistently uncovers connections both within and across layers under near-optimal network sparsity conditions. They also establish a consensus clustering of observations by applying a $k$-means algorithm to the local membership decomposition matrix.

## 3 Mixture of Multilayer SBM

Our model builds on SBM and considers two sets of latent variables corresponding respectively to the structure of the observations and the structure of the views. This proposal is at the crossroads of MLSBM, which involves the discovery of a traversing membership matrix of observations spanning all layers, and MMLSBM, which involves uncovering structural patterns within the layers.

**Observations.** We consider the observed data to be a tensor $\mathbf{A} \in \{0,1\}^{N \times N \times V}$ where $N$ is the number of observations (vertices), and $V$ the number of views. Each of the $V$ slices of $\mathbf{A}$ is an adjacency matrix corresponding to a graph $\mathcal{G}^v$. The tensor is thus a stack of adjacency matrices for multiple view graphs $(\mathcal{G}^1, \cdots, \mathcal{G}^V)$ with corresponding vertices. Let denote $(i,j)$ an edge between observations $i$ and $j$, we have by definition $A_{ijv} = \mathbb{I}_{((i,j) \in E^v)}$ where $E^v$ is the set of edges of the graph $\mathcal{G}^v$.

**Latent structures.** Let $\mathbf{Z} \in \{0,1\}^{N \times K}$ be the indicator membership matrix of observations, where $K$ is the number of view traversing clusters. We have by definition $Z_{ik} = \mathbb{I}_{(i \in k)}$, where $i$ denotes an observation and $k$ is a cluster across the views.

Let denote $\mathbf{W} \in \{0,1\}^{V \times Q}$, the indicator membership matrix for the views where $Q$ is the number of components of the view mixture. We have $W_{vs} = \mathbb{I}_{(v \in s)}$, where $v$ is a view and $s$ a cluster of views.

### 3.1 A mixture of observations through a mixture of views

The $V$ views are assumed to be generated by a mixture model of $Q$ components. Each component $s$ is a SBM. Each line of matrix $\mathbf{W}$ is assumed to follow a multinomial distribution, $\mathbf{W}_v \sim \mathcal{M}(1, \boldsymbol{\rho} = (\rho_1, \dots, \rho_Q))$, with

$$\mathbb{P}(\mathbf{W} \mid \boldsymbol{\rho}) = \prod_{v=1}^{V} \prod_{s=1}^{Q} \rho_v^{W_{vs}}.$$

Although we use multiple views with their known cluster structure (SBM), we assume a *traversing structure across all views* described by the latent variable $\mathbf{Z}$. By leveraging all available sources of information, we aim to achieve a more comprehensive understanding of the data and obtain community structures that are consistent across all views. The individuals are thus assumed to come from a number $K$ of sub-populations.

Each latent class vector for the observation follows a multinomial distribution, with $\mathbf{Z}_i \sim \mathcal{M}(1, \boldsymbol{\pi} = (\pi_1, \dots, \pi_K))$, and

$$\mathbb{P}(\mathbf{Z} \mid \boldsymbol{\pi}) = \prod_{i=1}^{N} \prod_{k=1}^{K} \pi_k^{Z_{ik}}.$$

Each observation $A_{ijv}$ conditionally to the latent structure $\mathbf{Z}$ follows a Bernoulli distribution: $A_{ijv} \mid Z_{ik} = 1, Z_{jl} = 1 \sim \mathcal{B}(\alpha_{kls})$. The probability of all observations given the latent variables $\mathbf{Z}$, $\mathbf{W}$ and a vector of parameters $\boldsymbol{\Theta}$, is thus

$$\mathbb{P}(\mathbf{A} \mid \mathbf{Z}, \mathbf{W}, \boldsymbol{\Theta}) = \prod_{\substack{i=1, \, k=1 \\ i<j \;\; l=1}}^{N \quad K} \prod_{v=1}^{V} \prod_{s=1}^{Q} \left( \alpha_{kls}^{A_{ijv}} (1 - \alpha_{kls})^{1-A_{ijv}} \right)^{Z_{ik} Z_{jl} W_{vs}}.$$

### 3.2 Identifiability

**Theorem 1** *Let $N \geq \max(2K, 4Q)$ and $V \geq 2K$. Assume that for any $k, l \in \{1, \dots, K\}$ and every $s \in \{1, \dots, Q\}$, the coordinates of $\boldsymbol{\pi}^T \boldsymbol{\alpha}_{k..} \boldsymbol{\rho}$ are all different, $(\boldsymbol{\pi}^T \boldsymbol{\alpha}_{..s} \boldsymbol{\pi})_{s=1:Q}$ are distinct, and each $(\boldsymbol{\alpha}_{kl.} \boldsymbol{\rho})_{k,l=1:K}$ differs. Then, the mimi-SBM parameters $\boldsymbol{\Theta} = (\boldsymbol{\pi}, \boldsymbol{\rho}, \boldsymbol{\alpha})$ are identifiable.*

**Proof** The proof of this theorem is given in Appendix A. ∎

### 3.3 Bayesian modeling

Bayesian modeling provides a natural framework for incorporating prior knowledge which can improve the accuracy of the estimated block structure, mainly when the available data is limited or noisy.

In this context, we follow Latouche et al. (2012) and define the chosen conjugate distributions both for the proportion of the mixture and the proportions of the blocks. Conjugate

priors lead to closed-form posterior distributions.

$$\mathbb{P}(\boldsymbol{\pi} \mid \boldsymbol{\beta}^0 = (\beta_1^0, \ldots, \beta_K^0)) = \text{Dir}(\boldsymbol{\pi}; \boldsymbol{\beta}^0), \tag{1}$$

$$\mathbb{P}(\boldsymbol{\rho} \mid \boldsymbol{\theta}^0 = (\theta_1^0, \ldots, \theta_Q^0)) = \text{Dir}(\boldsymbol{\rho}; \boldsymbol{\theta}^0), \tag{2}$$

where Dir(.) stands for the Dirichlet distribution.

$$\mathbb{P}(\boldsymbol{\alpha} \mid \boldsymbol{\eta}^0 = (\eta_{kls}^0), \boldsymbol{\xi}^0 = (\xi_{kls}^0)) = \prod_{k,k<l} \prod_s \text{Beta}(\alpha_{kls}; \eta_{kls}^0, \xi_{kls}^0). \tag{3}$$
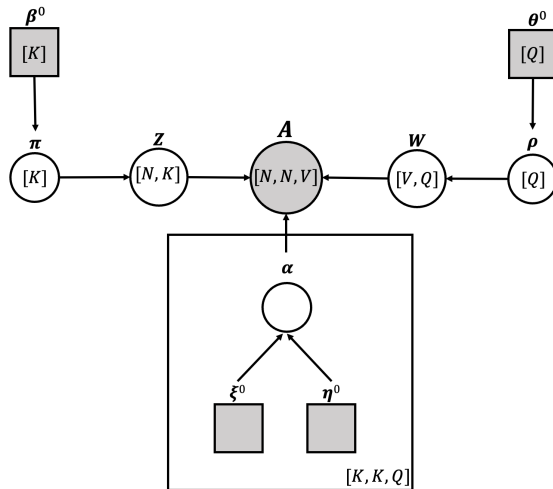


Figure 2: Illustration of mimi-SBM with bayesian notations

The parameters $\boldsymbol{\beta}^0, \boldsymbol{\theta}^0, \boldsymbol{\eta}^0, \boldsymbol{\xi}^0$ are chosen according to Jeffreys priors which are often considered non-informative or weakly informative. They do not introduce strong prior assumptions or biases into the analysis.

For the Dirichlet distribution, a suitable choice for $\beta_k^0$ and $\theta_s^0$ is setting them both to 1/2, which directly corresponds to an objective Jeffreys prior distribution. Similarly, for the Beta distribution, $\eta_{kls}^0$ and $\xi_{kls}^0$ can be chosen as 1/2 for all appropriate indices k, l, and s.

## 4 Variational Bayes Expectation Maximisation for mimi-SBM

Computing the marginal likelihood is a challenging problem in Stochastic Block Models:

$$\mathbb{P}(\mathbf{A}) = \sum_{\mathbf{Z}} \sum_{\mathbf{W}} \int \int \int \mathbb{P}(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho}) \, d\boldsymbol{\alpha} \, d\boldsymbol{\pi} \, d\boldsymbol{\rho}. \tag{4}$$

The computation of integrals in the formula for this marginal likelihood presents analytical challenges or becomes unfeasible, while sums over $\mathbf{Z}$ and $\mathbf{W}$ become impractical when the number of parameters or observations is substantial.

Approximation of complex posterior distributions is usually performed either by sampling (Markov Chain Monte Carlo or related approaches) or by Variational Bayes inference

introduced by Attias (1999). Variational Bayes Expectation Maximization algorithm offers several advantages such as reduced computation time and the ability to work on larger databases.

## 4.1 Evidence Lower Bound

Variational inference is computationally efficient and scalable to large datasets and work especially well for SBM models. It formulates the problem as an optimization task, where the goal is to find the best approximation to the true posterior distribution. This optimization framework allows for efficient computation of the variational parameters by maximizing a lower bound on the log-likelihood, known as the evidence lower bound (ELBO). Optimization techniques like stochastic gradient descent (SGD) or Expectation-Maximisation algorithm can be employed to find the optimal variational parameters.

The distribution $\mathbb{P}(\mathbf{Z}, \mathbf{W} | \mathbf{A}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})$ is intractable when taking SBM into account, hence we approximate the entire distribution $\mathbb{P}(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho} \mid \mathbf{A})$. Given a variational distribution $q$ over $\{\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho}\}$, we can decompose the marginal log-likelihood into Evidence Lower BOund (ELBO) part and KL-divergence between variational and posterior distribution :

$$
\begin{aligned}
\log P(\mathbf{A}) &= \mathbb{E}_q \left[ \log \frac{P(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})}{P(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho} | \mathbf{A})} \right] \\
&= \underbrace{\mathbb{E}_q \left[ \log \frac{P(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})}{q(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})} \right]}_{ELBO = \mathcal{L}(q(.))} + \underbrace{\mathbb{E}_q \left[ \log \frac{q(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})}{P(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho} \mid \mathbf{A})} \right]}_{\mathbf{KL}(q(\cdot) \| \mathbb{P}(\cdot | \mathbf{A}))}
\end{aligned}
$$

where $\mathbf{KL}\left(q(.) \mid \mathbb{P}(\cdot \mid \mathbf{A})\right) = -\mathbb{E}_q[\log \frac{p}{q}] \geq -\log \mathbb{E}_q[\frac{p}{q}] \geq 0$ from Jensen inequality.

The ELBO is given by

$$
\mathcal{L}(q(\cdot)) = \sum_{\mathbf{Z}, \mathbf{W}} \int \int \int q(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho}) \log \frac{p(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})}{q(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})} \, d\boldsymbol{\alpha} \, d\boldsymbol{\pi} \, d\boldsymbol{\rho}. \tag{5}
$$

The variational distribution is typically selected from an easier-to-handle family of distributions, such as the exponential family. The variational distribution's parameters are then adjusted to reduce the KL divergence to the posterior distribution. If $q(.)$ is exactly $p(.|\mathbf{A})$ the $\mathbf{KL}$ term is equal to 0, and the ELBO is maximized.

We assume a mean-field approximation for $q(\cdot)$:

$$
\begin{aligned}
q(\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho}) &= \prod_{i=1}^{N} q(\mathbf{Z}_i) \prod_{v=1}^{V} q(\mathbf{W}_v) \prod_{s=1}^{Q} \prod_{k, k \leq l}^{K} q(\alpha_{kls}) \, q(\boldsymbol{\pi}) \, q(\boldsymbol{\rho}) \\
&= \operatorname{Dir}(\boldsymbol{\pi}; \boldsymbol{\beta}) \, \operatorname{Dir}(\boldsymbol{\rho}; \boldsymbol{\theta}) \prod_{i=1}^{N} \mathcal{M}(\mathbf{Z}_i; 1, \boldsymbol{\tau}_i) \prod_{v=1}^{V} \mathcal{M}(\mathbf{W}_i; 1, \boldsymbol{\nu}_v) \\
&\qquad \prod_{s=1}^{Q} \prod_{k, k \leq l}^{K} \operatorname{Beta}(\alpha_{kls}; \eta_{kls}, \xi_{kls}),
\end{aligned} \tag{6}
$$

where $\tau_{ik}$ (resp. $\nu_{vs}$) are variational parameters indicating the probability that individual $i$ (resp. a view $v$) belongs to cluster $k$ (resp. component $s$).

9

According to (5), given a distribution $q(.)$, the ELBO is given by

$$
\begin{aligned}
\mathcal{L}\left(q(.)\right) = & \log\left\{\frac{\Gamma\left(\sum_{k=1}^{K}\beta_k^0\right)\prod_{k=1}^{K}\Gamma\left(\beta_k\right)}{\Gamma\left(\sum_{k=1}^{K}\beta_k\right)\prod_{k=1}^{K}\Gamma\left(\beta_k^0\right)}\right\} + \log\left\{\frac{\Gamma\left(\sum_{s=1}^{Q}\theta_s^0\right)\prod_{s=1}^{Q}\Gamma\left(\theta_s\right)}{\Gamma\left(\sum_{s=1}^{Q}\theta_s\right)\prod_{s=1}^{Q}\Gamma\left(\theta_s^0\right)}\right\} \\
& + \sum_{k\leq l}^{K}\sum_{s=1}^{Q}\log\left\{\frac{\Gamma\left(\eta_{kls}^0 + \xi_{kls}^0\right)\Gamma\left(\eta_{kls}\right)\Gamma\left(\xi_{kls}\right)}{\Gamma\left(\eta_{kls} + \xi_{kls}\right)\Gamma\left(\eta_{kls}^0\right)\Gamma\left(\xi_{kls}^0\right)}\right\} \\
& - \sum_{i}^{N}\sum_{k}^{K}\tau_{ik}\log\tau_{ik} - \sum_{v}^{V}\sum_{s}^{Q}\nu_{vs}\log\nu_{vs},
\end{aligned}
\tag{7}
$$

where $\Gamma(.)$ is the Gamma function. This function is also called Integrated Likelihood variational Bayes (ILvb, Latouche et al. (2012)), since it can be used for model selection.

### 4.2 Lower bound optimization

We consider a Variational Bayes EM algorithm for estimating the parameters. The algorithm starts by initializing the model parameters and then iteratively performs two steps: the Variational Bayes Expectation step (VBE-step) and the Maximization step (M-step) (See Algorithm 1).

In the VBE-step, the variational distributions are optimized over latent variables: $q(\mathbf{Z}_i)\ \forall i \in \{1,\dots,N\}$ and $q(\mathbf{W}_v)\ \forall v \in \{1,\dots,V\}$ to approximate the true posterior distribution.

In the M-step, the parameters of the model are updated to maximize a lower bound on the log-likelihood, with respect to parameters computed in the VBE-step: $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, $\boldsymbol{\eta}$, and $\boldsymbol{\xi}$.

There exist multiple techniques for initializing the EM algorithm. One prevalent approach involves using random initial values, where the model parameters are assigned random values drawn from a designated distribution. Nevertheless, this method may lack reliability and fail to provide satisfactory starting values for the algorithm. According to the initialization method proposed in Stanley et al. (2016), the parameters $(\tau_{ik})$ and $(\nu_{vs})$ are initialized based on the outcomes of a stochastic block model (SBM) applied separately to each view. The objective is to capture the overall structure of the data from each view using SBM, combine this information using K-means clustering, and subsequently refine the obtained results using our model.

## 5 Model selection

In the context of clustering, model selection often refers to the process of determining the ideal number of clusters for a given dataset. In our situation, the key decision lies in selecting appropriate values for $K$ and $Q$ to strike a balance between data attachment and model complexity. To achieve this, several criteria based on penalized log-likelihood can be employed, such as the Akaike Information Criterion (AIC) (Akaike, 1998), Bayesian Information Criterion (BIC) (Schwarz, 1978) and more recently the Integrated Completed Likelihood (ICL) (Biernacki et al., 2000). We specifically consider the ICL criterion and its associated penalties as they frequently yield good trade-offs in the selection of mixture models (Biernacki et al., 2010).

---

**Algorithm 1** mimi-SBM

---

**Require:** Tensor of adjacency matrices $\mathbf{A}$, Number of clusters $K$, Number of components of the views $Q$, precision *eps*.

Initialization : $\tau_{ik}^{(old)}$ and $\nu_{ik}^{(old)}$

**while** $\| \mathcal{L}\left(q^{new}(.)\right) - \mathcal{L}\left(q^{old}(.)\right) \| < eps$ **do**

   **VBE-step**

   Compute $\tau_{ik}^{(new)}$ $\forall i \in \{1, \ldots, N\}$ and $\forall k \in \{1, \ldots, K\}$

   Compute $\nu_{vs}^{(new)}$ $\forall v \in \{1, \ldots, V\}$ and $\forall s \in \{1, \ldots, Q\}$

   **M-step**

   Optimize $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, $\boldsymbol{\eta}$, $\boldsymbol{\xi}$ with respect to $(\tau_{ik}^{(new)})$ and $(\nu_{vq}^{(new)})$

   **ELBO**

   Compute $\mathcal{L}\left(q^{new}(.)\right)$

**end while**

---

The ICL is based on the log-likelihood integrated over the parameters of the complete data. Furthermore, if we assume that parameters of component-connection probability $\boldsymbol{\alpha}$, parameters for mixture of communities $\boldsymbol{\pi}$ and parameters for views mixture $\boldsymbol{\rho}$ are independent, we have:

$$
\begin{aligned}
\text{ICL}(\mathbf{A}, K, Q) &= \log \mathbb{P}(\mathbf{A}, \mathbf{Z}, \mathbf{W} \mid K, Q) \\
&= \log \int_{\boldsymbol{\alpha}} \mathbb{P}(\mathbf{A} \mid \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}) \, \mathbb{P}(\boldsymbol{\alpha}) d\boldsymbol{\alpha} \\
&\quad + \log \int_{\boldsymbol{\pi}} \mathbb{P}(\mathbf{Z} \mid \boldsymbol{\pi}) \, \mathbb{P}(\boldsymbol{\pi}) d\boldsymbol{\pi} \\
&\quad + \log \int_{\boldsymbol{\rho}} \mathbb{P}(\mathbf{W} \mid \boldsymbol{\rho}) \, \mathbb{P}(\boldsymbol{\rho}) d\boldsymbol{\rho}.
\end{aligned}
\tag{8}
$$

In our variational framework, $\mathbf{Z}$ and $\mathbf{W}$ must be estimated. $\hat{\mathbf{Z}}$ (resp. $\hat{\mathbf{W}}$) can be chosen as the variational parameters $\boldsymbol{\tau}$ (resp. $\boldsymbol{\nu}$) directly or by a Maximum a Posteriori (MAP):

$$
\hat{\mathbf{Z}}_i = \underset{k \in 1:K}{\operatorname{argmax}} \ \tau_{ik}.
$$

By using approximations, such as Stirling's approximation formula on $\mathbb{P}(\boldsymbol{\pi})$ and $\mathbb{P}(\boldsymbol{\rho})$ and the Laplace asymptotic approximation on $\mathbb{P}(\boldsymbol{\alpha})$, we can define an *approximate ICL*:

$$
\begin{aligned}
\text{ICL}(\mathbf{A}, K, Q) &\approx \log \mathbb{P}(\mathbf{A}, \hat{\mathbf{Z}}, \hat{\mathbf{W}} \mid K, Q) - \text{pen}(K, Q) \\
&\approx \mathcal{L}\left(q(.)\right) - \text{pen}(K, Q),
\end{aligned}
\tag{9}
$$

where

$$
\text{pen}(K, Q) = \frac{1}{2} \frac{K(K+1)}{2} Q \log(V \frac{N(N-1)}{2}) + \frac{1}{2}(K-1)\log(N) + \frac{1}{2}(Q-1)\log(V).
$$

The penalization in this *approximate ICL* is composed of a part depending on the number of parameters of component-connection probability tensor $\boldsymbol{\alpha}$ and the number of vertices taken into account, and a part that takes into account the number of degree of freedom in

mixture parameters and the number of variables related to them. Recall that our model is based on undirected (symmetric) adjacency matrices, so we only consider the upper triangular matrices (without the diagonal).

However, in the Bayesian framework with conjugate priors, it is possible to define an exact ICL (Côme and Latouche, 2015). Moreover, it can be obtained from the previously defined ILvb (7) when the entropy of the latent variables is zero and the Expectation-Maximization algorithm is a Classification EM (CEM, Celeux and Govaert, 1992). In other words, variational parameters are equal to 1 if it is the MAP and 0 otherwise. Thus, this *exact ICL* can be defined as:

$$
\begin{aligned}
\mathrm{ICL}_{\mathrm{exact}}(\mathbf{A}, K, Q) = & \log \left\{ \frac{\Gamma\left(\sum_{k=1}^{K} \beta_k^0\right) \prod_{k=1}^{K} \Gamma(\beta_k)}{\Gamma\left(\sum_{k=1}^{K} \beta_k\right) \prod_{k=1}^{K} \Gamma(\beta_k^0)} \right\} \\
& + \log \left\{ \frac{\Gamma\left(\sum_{s=1}^{Q} \theta_s^0\right) \prod_{s=1}^{Q} \Gamma(\theta_s)}{\Gamma\left(\sum_{s=1}^{Q} \theta_s\right) \prod_{s=1}^{Q} \Gamma(\theta_s^0)} \right\} \\
& + \sum_{k \leq l}^{K} \sum_{s=1}^{Q} \log \left\{ \frac{\Gamma\left(\eta_{kls}^0 + \xi_{kls}^0\right) \Gamma(\eta_{kls}) \Gamma(\xi_{kls})}{\Gamma\left(\eta_{kls} + \xi_{kls}\right) \Gamma(\eta_{kls}^0) \Gamma(\xi_{kls}^0)} \right\}.
\end{aligned}
\tag{10}
$$

Instead of using a CEM, it is possible to use the v Variational parameters directly, and to derive a *variational ICL* from the previous criterion. Figure 3 summarizes the links between the various selection criteria and clearly shows that in a particular context *exact ICL*, *Variational ICL*, and *ILvb* criteria are identical.
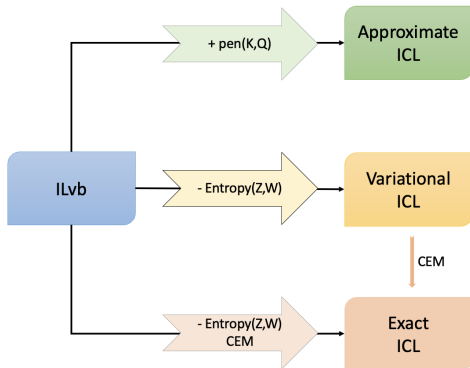


Figure 3: Diagram of links between different model selection criteria.

## 6 Experiments

### 6.1 Simulation study

To ensure that *mimi-SBM* behaves consistently, we have developed a simulation scheme, as depicted in Figure 4. The simulated data have been designed to reflect the complexity found in real-world data clustering challenges. In particular, this scheme allows for diverse

clustering patterns across different view components. Also, it includes the possibility of controlling clustering errors, with observations being inaccurately assigned to an incorrect group, implying inconsistencies in the adjacency matrices.

**Simulated data.** Artificial adjacency matrices are generated from observations and views. We aim to establish a link between the simulated adjacency matrices and the final clustering that most accurately represents a problem of meta-clustering.

Various parameter values are tested for $N$ (observations), $V$ (views), $K$ (clusters), and $Q$ (components) in different scenarios. Besides, $\boldsymbol{\pi}, \boldsymbol{\rho}$ correspond to an equiprobability of belonging to a cluster or component, thus $\{\pi_k\}_{k=1}^K = 1/K$ and $\{\rho_s\}_{s=1}^Q = 1/Q$ (Figure 4, *Parameters*).

First of all, it is assumed that the number of real clusters ($K$) will always be equal or higher than the number of clusters coming from each component. Also, each component has a precise number of clusters ($K^q$), and each view belonging to this component will have this number of clusters $K^q \sim \mathcal{U}(\{2, \ldots, K\})$, where $\mathcal{U}$ is the discrete uniform distribution (Figure 4, *Clusters per component*).

Now, for each component, we randomly associate a link between the final consensus clusters $\mathbf{Z}$ and clusters coming from the component $\mathbf{Z}^q$, ensuring that no component cluster remains empty (Figure 4, *Links between clusters*).

Eventually, for each pair of nodes $(i, j)$ and each layer $v$, an edge is generated with probability $\alpha_{Z_i Z_j W_v}$, leading to set the corresponding entry in the multilayer adjacency tensor $\mathbf{A}$ to 1 or 0 (Figure 4, *Generation of edges*).

The simulated data are used to assess three aspects:

1. **Model selection.** In Section 6.1.1, various criteria are examined to to recover the true parameters $K$ and $Q$.

2. **Clustering ability.** In Sections 6.1.2 and 6.1.3, *mimi-SBM* is evaluated against other state-of-the-art techniques regarding the clustering of observations and the clustering of views using ARI scores (see below).

3. **Robustness.** In Section 6.1.4, we investigate further the model ability to handle noisy configurations inherent in real-world clustering problems.

The code for the simulations is available on the CRAN, and on GitHub in the repository *mimiSBM*. [1]

**Adjusted Rand Index.** The Adjusted Rand Index (ARI, Hubert and Arabie, 1985) quantifies the similarity between two partitions. In the simulation to follow, it quantifies the similarity between the prediction by our clustering models and the true partition. It corresponds to the proportion of pairs $(i, j)$ of observations jointly grouped or separated. The more similar the partitions, the closer the ARI is to 1.

6.1.1 Comparing model selection criteria

In this section, our goal is to undertake a comparative analysis of model selection criteria to determine the optimal choice of criterion.
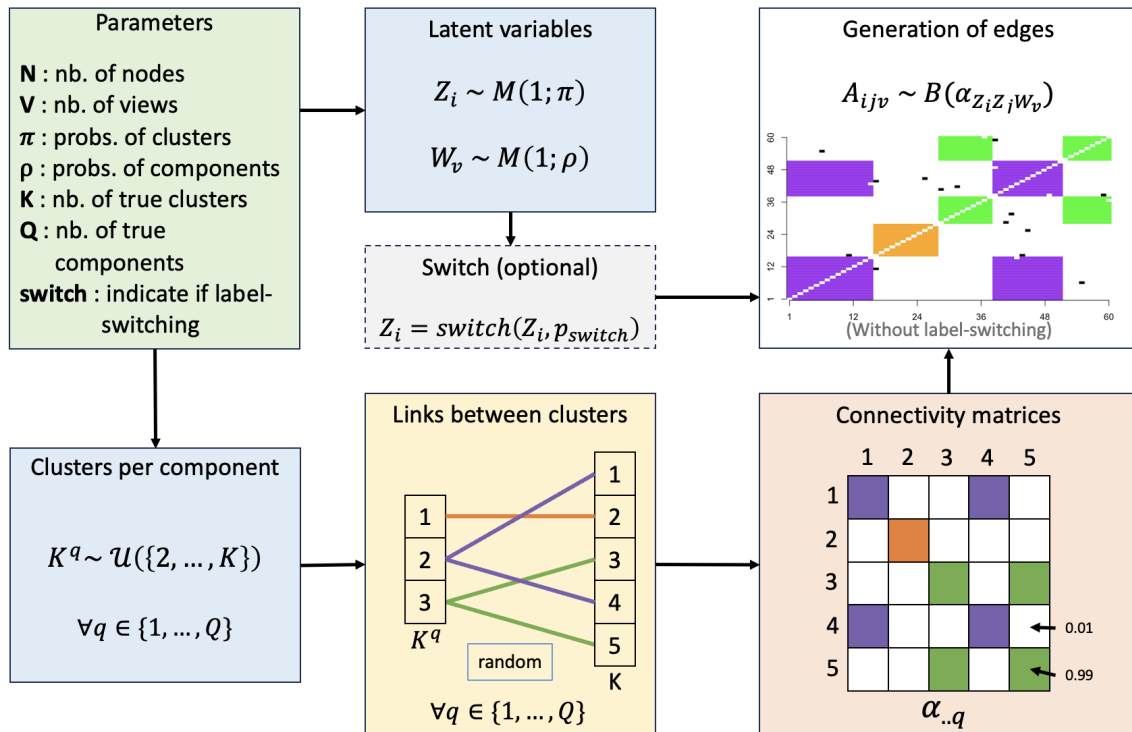
---

1. https://github.com/Kdesantiago/mimiSBM.

Figure 4: Diagram of the simulation process. Example of adjacency matrices resulting from mixing and traversing clusters across views, with potentially label-switching, for $K = 5$. For each view component, a number of clusters $K^q$ is randomly drawn (discrete uniform distribution). Each cluster in the $q^{\text{th}}$ component is then linked to certain clusters in the final partition. For this component, $K^q = 3$, and final clusters 1 (respectively 3) and 4 (resp. 5) are merged into cluster 2 (resp. 3) of the component, and the first cluster of the component corresponds perfectly to the final consensus cluster 2. Afterwards, these links are represented by a very strong connectivity within the $\boldsymbol{\alpha}_{..q}$ matrix ($p = 0.99$) and a very weak one ($p = 0.01$) for the others.

The use of simulations gives us a complete control over the hyperparameters that generated the data. To do this, we generated 50 different datasets with hyperparameters $K = 10$ and $Q = 5$. The model selected for each criterion is the one that maximizes its value.

Simulation results in Figure 5, clearly show that, in all scenarios, each criterion delivers consistent and comparable performance. Without exception, the criteria consistently produce the same selection of clusters and number of view components.

In Figure 5a, only the hyperparameter $K$ varies. This parameter was mostly well estimated because, during model selection, the true number of clusters was typically identified in the majority of cases. However, it was crucial to note that in the context of individual clustering, the criteria tended to overestimate the number of clusters.

In Figure 5b, the number of components parameter $Q$ is variable, while $K$ remains constant. In the majority of scenarios, it was observed that the number of components was accurately estimated. Furthermore, when a fixed parameter for clustering was considered,

(a) Model selection on $K$, with $Q$ fixed



(b) Model selection on $Q$, with $K$ fixed



(c) Model selection on $K$, with $Q$ free



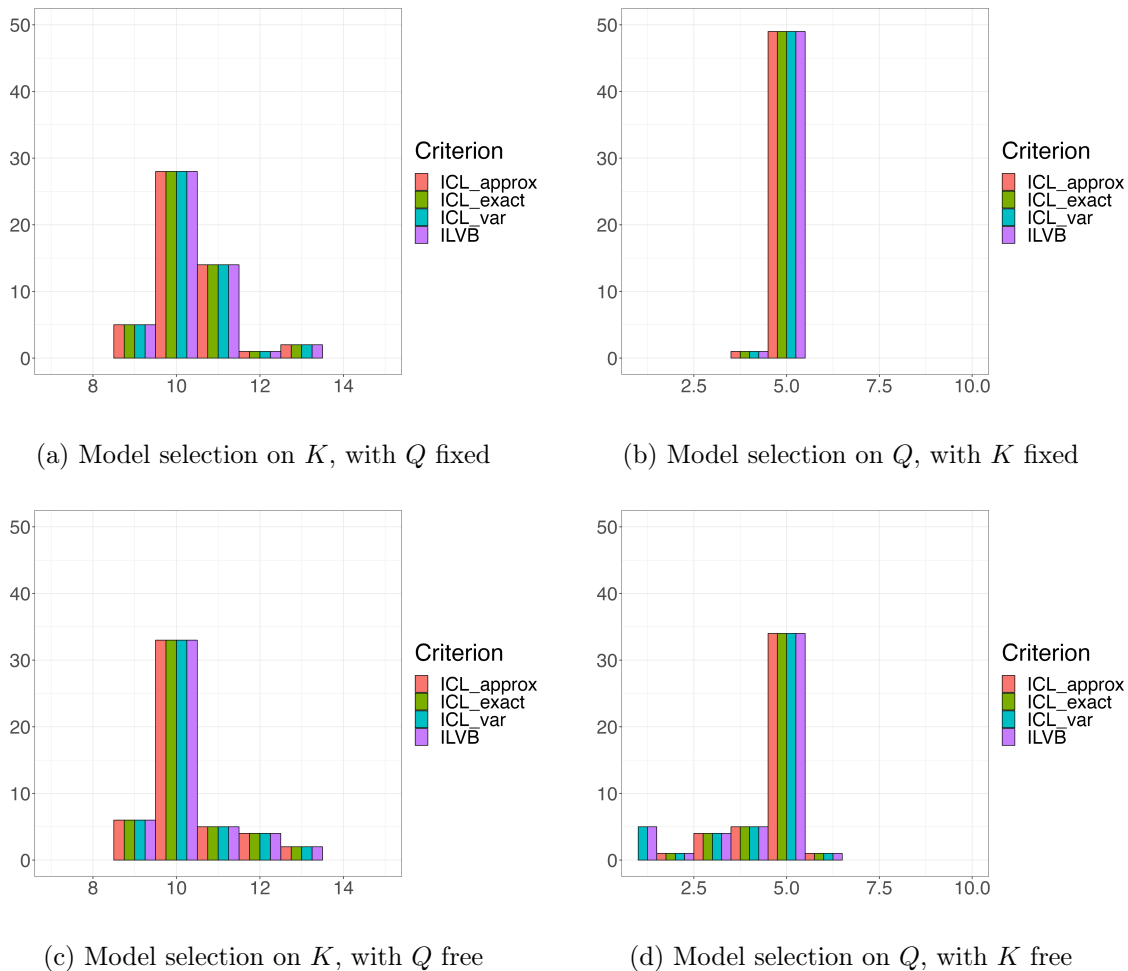(d) Model selection on $Q$, with $K$ free

Figure 5: Bar plots of model selection criteria on 50 simulations with 10 true clusters for observations and 5 true components for views. Figure (a) (resp. (b)) indicates the number of times the $K$ (resp. $Q$) value selected while the other parameters is set to the true value. Figures (c) and (d) show the same information when hyperparameters are optimized at the same time.

the task inherently became more tractable due to the use of abundant information for the estimation of view components.

In Figures 5c and 5d, the selection of hyperparameters is aligned with fixed-parameter results. The criteria consistently demonstrate an aptitude for identifying the optimal cluster and view component quantities. Nonetheless, akin to previous instances, the model occasionally exhibits errors in hyperparameter estimation, often underestimating the number of components while overestimating the number of classes.

6.1.2 SIMULATIONS WITHOUT LABEL-SWITCHING

**Comparison of clustering.** In Figure 6, it has been observed that *mini-SBM* achieved the best clustering results for each considered experimental configuration. Indeed, *mini-SBM* recorded the highest ARI score for all data sizes, number of clusters and sources. On the other hand, the *M3C* model improved as the number of views, clusters and sources increased. In contrast, the *TWIST* model showed poorer performances as the clustering problem became more complex, suggesting that this model may be less suitable for difficult clustering problems.
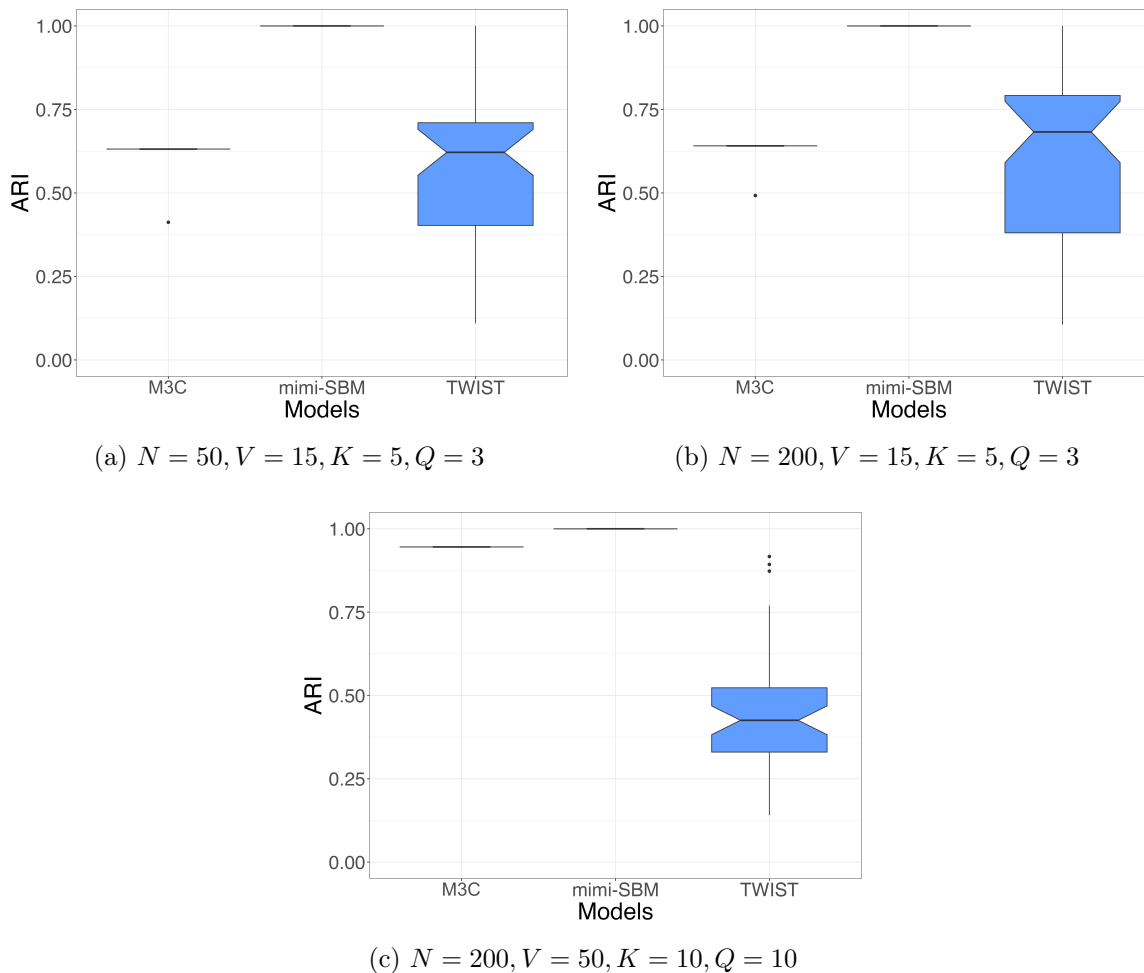


(a) $N = 50, V = 15, K = 5, Q = 3$

(b) $N = 200, V = 15, K = 5, Q = 3$

(c) $N = 200, V = 50, K = 10, Q = 10$

Figure 6: Boxplot of ARI measure between true partition and output partition of *M3C*, *mimi-SBM* and *TWIST* models.

**Comparison of view components.** In Figure 7, as the number of observations increases, the performances of the models generally tends to improve. However, the *graphclust* model appears to identify the true sources less frequently than the *mimi-SBM* and *TWIST* models.

While *TWIST* often identifies the true members of the sources perfectly, it does make some errors, visible as outliers on the boxplot.
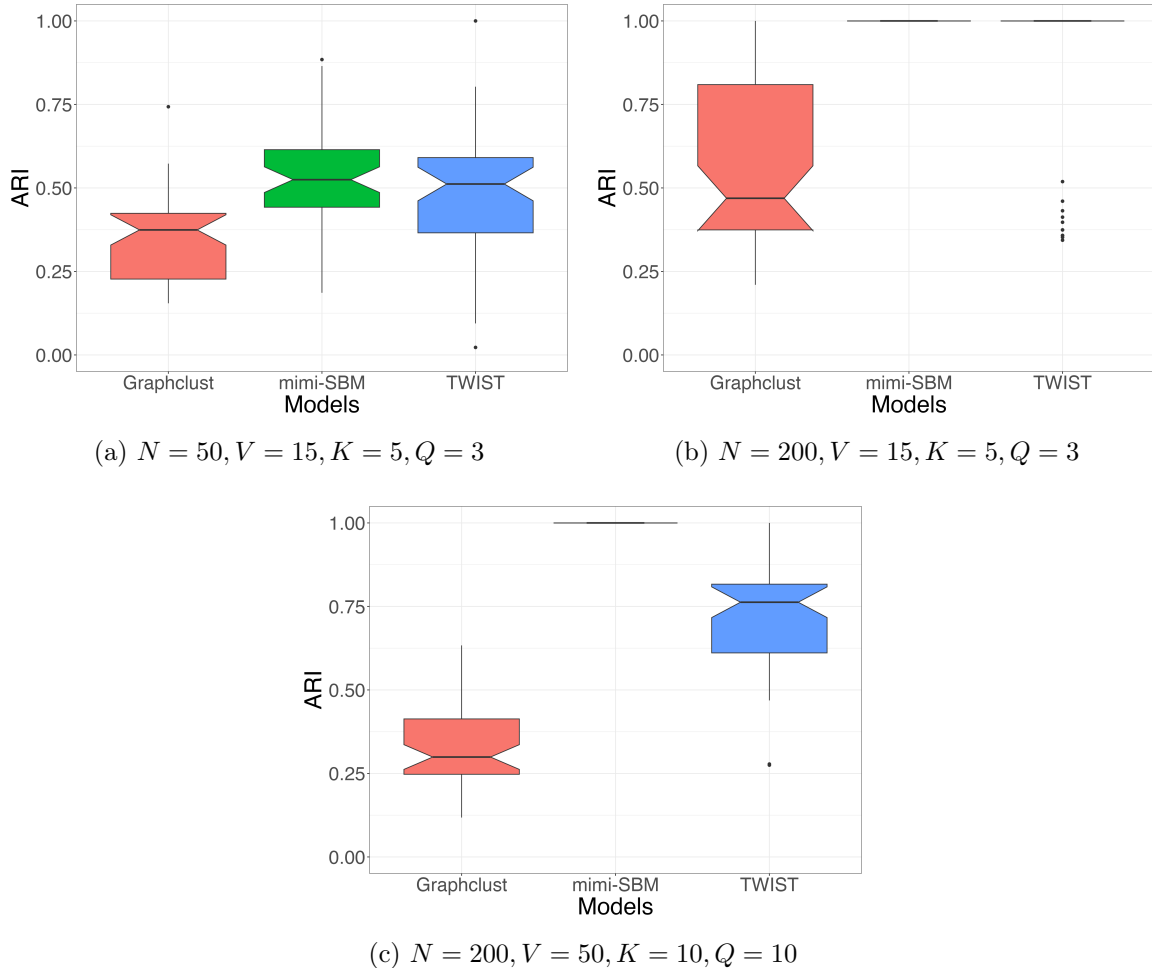


(a) $N = 50, V = 15, K = 5, Q = 3$       (b) $N = 200, V = 15, K = 5, Q = 3$

(c) $N = 200, V = 50, K = 10, Q = 10$

Figure 7: Boxplot of ARI measure between true view clustering and output clustering of *graphclust*, *mimi-SBM* and *TWIST* models.

### 6.1.3 SIMULATIONS WITH LABEL-SWITCHING

In this section, we revisit the analyses from the previous section, but with a focus on a more challenging issue: label-switching. The idea of perturbing the cluster labels within the generation process simulates the fact that an individual has been associated with another cluster during the process of creating adjacency matrices.

In our context, we simulate the fact that an individual belongs to the real cluster, and then we simulate the representation of this clustering by the link between the final clustering and the one specific to each view component, to obtain the different affinity matrices.

The perturbation occurs during the generation of individual component-based clusters. For each view, a perturbation is introduced for each individual with a probability of $p_{\text{switch}} =$

17

0.1. In such perturbation, the respective individual is then associated with one of the other available clusters. As a result, the probability of creating a link between individuals is influenced.

**Comparison of clustering.** The analysis summarized in Figure 8 reveals that across all examined experimental setups, the *mimi-SBM* consistently attained the most favorable clustering outcomes. Furthermore, it is noteworthy that the score variability associated with *mimi-SBM* is notably lower than that observed for other models. The effectiveness of the *M3C* and *TWIST* model showed improvement as the number of views, clusters, and sources increased, yet it maintained a relatively high level of variance. The number of individuals to be clustered plays a crucial role in minimizing errors. This effect stems from the fact that a larger number of individuals subject to clustering contributes to a more robust estimation of the parameters.
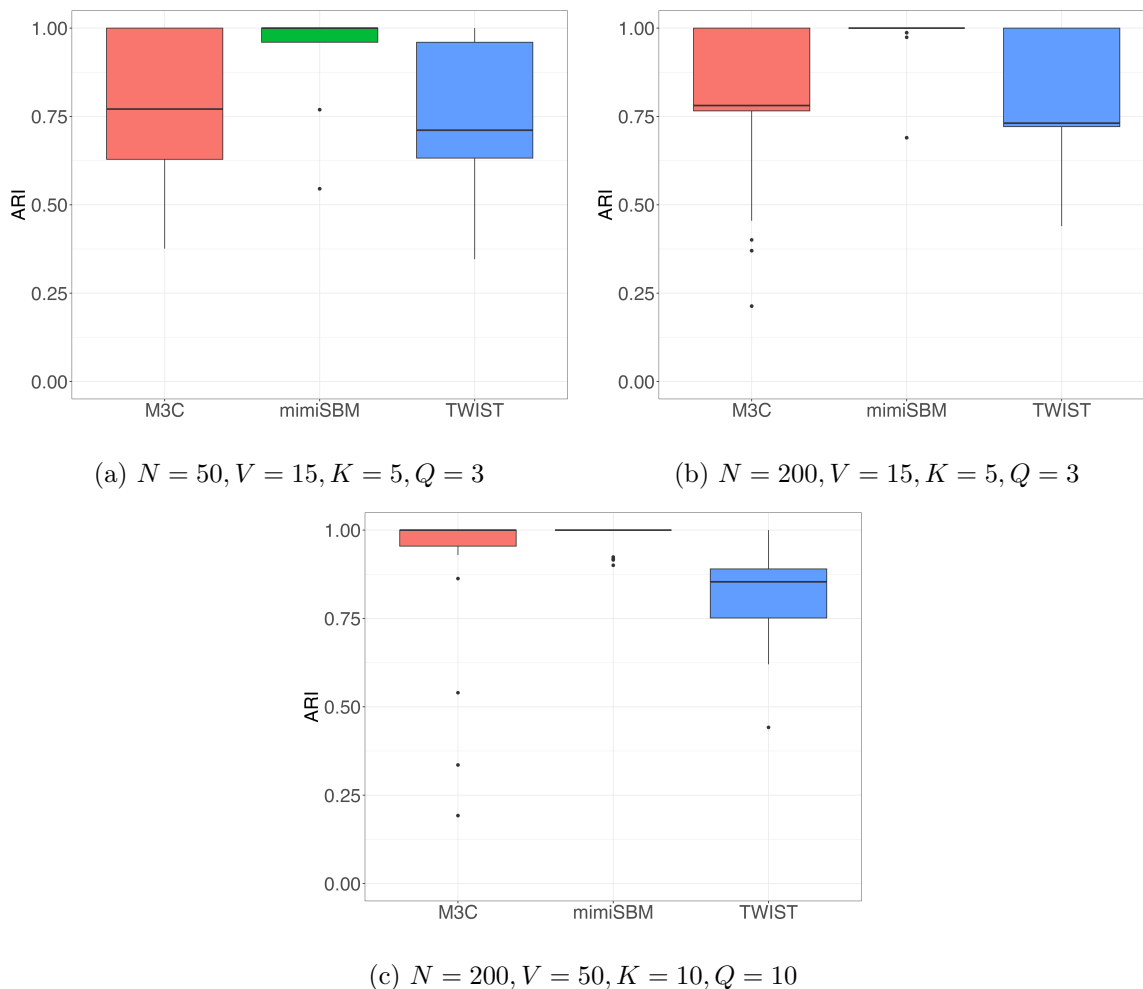
(a) $N = 50, V = 15, K = 5, Q = 3$

(b) $N = 200, V = 15, K = 5, Q = 3$

(c) $N = 200, V = 50, K = 10, Q = 10$

Figure 8: Boxplot of ARI measure between true partition and output partition of *M3C*, *mimi-SBM* and *TWIST* models.

**Comparison of view components.** In Figure 9, as the quantity of observations increases, the models typically exhibit enhanced performance. Nevertheless, the *graphclust* model seems to less frequently pinpoint the actual sources compared to the *mimi-SBM* and *TWIST* models. When faced with a small number of perspectives, *TWIST* model displays significant variability. While it consistently delivers good results, it remains vulnerable to unfavorable initializations, which can lead to notably suboptimal clustering outcomes. Moreover, when label-switching is introduced, the model's performance is observed to be slightly less effective compared to the precedent scenario.

Similar to the scenario without label-switching, the model experiences considerable variability in its estimation when dealing with a limited number of individuals and perspectives. However, as the number of individuals and views increases, the variance of ARI decreases noticeably, accompanied by an improvement in performance. *Mimi-SBM* model consistently demonstrates efficacy across all cases, even including perturbations in the adjacency matrices used for clustering.

### 6.1.4 ROBUSTNESS TO LABEL-SWITCHING

Given that the *mimi-SBM* showed a satisfactory performance level in the previous section, even under the influence of label-switching perturbation, this section aims to further assess the robustness and limitations of our model concerning this criterion.

By varying the label switching rate, from 0 to 1 in steps of 0.10, in order to see the evolution of clustering capacities on individuals and views.

For Figures 10 and 11, clustering performances demonstrate a significant level of efficacy when the label-switching rate is low.

As the rate of switched labels exceeds 40%, the stability of the individual clustering process progressively diminishes. This trend continues until the clustering process becomes entirely arbitrary when the switch-labeling rate surpasses 60%, as contrasted with the true partition. An observable improvement in performance becomes evident as the switched label rate approaches 1. This outcome is logically anticipated, as the reassignment of all individuals from one cluster to another results in their distribution across $K-1$ clusters instead of the initial $K$ clusters.

In the context of view-based clustering, we encounter a similar set of observations, albeit with a much more pronounced decline in performance. When the label-switching rate surpasses 20%, the ability of *mimi-SBM* to effectively identify view components experiences a drastic reduction. Furthermore, when this rate exceeds 40%, the feasibility and relevance of conducting clustering based on these views are severely compromised. One plausible explanation for this phenomenon is that, due to the perturbation, each adjacency matrix becomes highly noisy, lacking any discernible structure. Consequently, the model struggles to distinguish any specific connections within the mixtures, leading to a notably diminished clustering performance score.

**Summary.** The *mimi-SBM* model has shown its capability in successfully recovering the stratification of individuals and the components of the mixture of views, even when the data is perturbed. However, like any statistical model, its performance, especially regarding the mixture of views, benefits from larger sample sizes. The accurate modeling of mixture

(a) $N = 50, V = 15, K = 5, Q = 3$

(b) $N = 200, V = 15, K = 5, Q = 3$

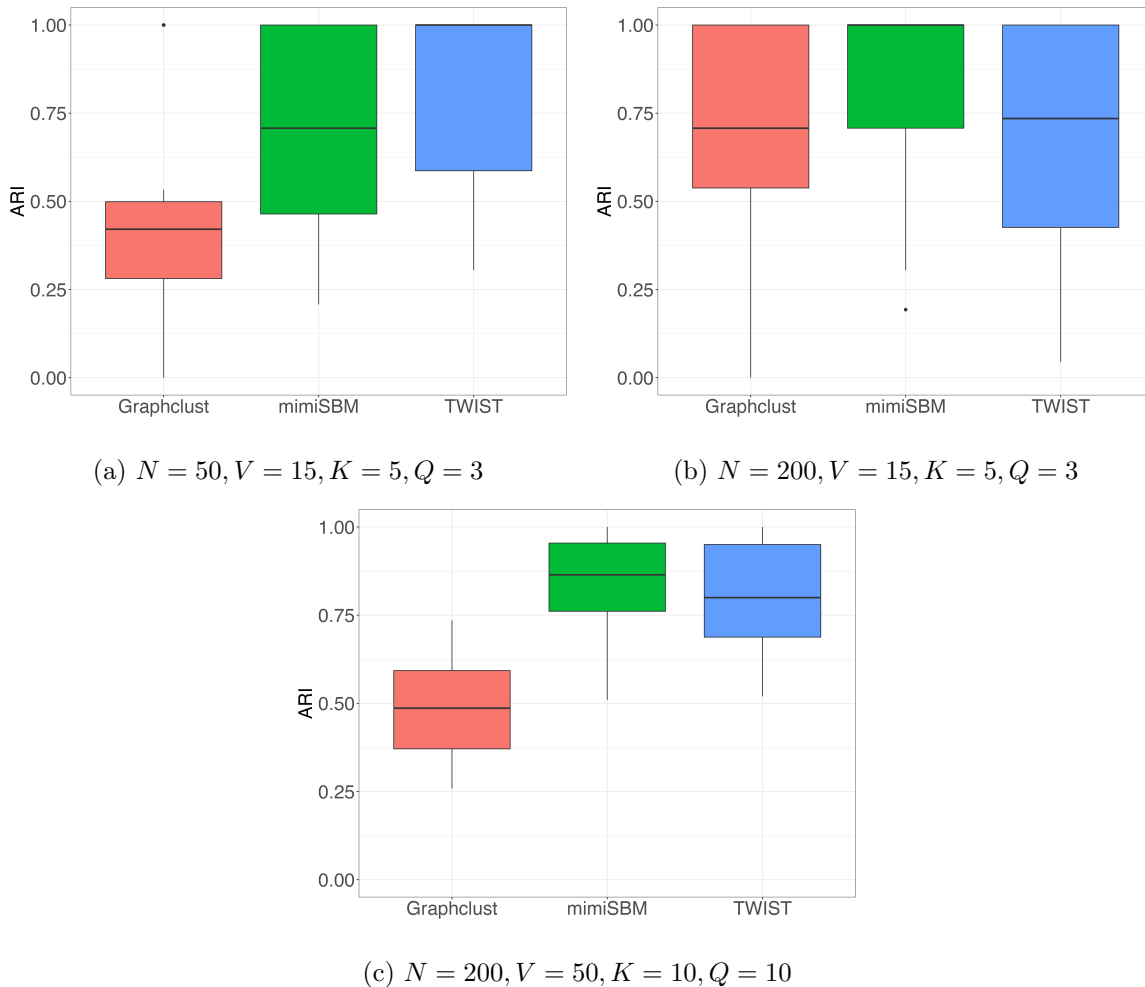(c) $N = 200, V = 50, K = 10, Q = 10$

Figure 9: Boxplot of ARI measure between true view clustering and output clustering of *graphclust*, *mimi-SBM* and *TWIST* models.

components is crucial in various applications, making the mimi-SBM model highly valuable in a wide range of contexts.

## 6.2 Worldwide Food Trading Networks

**Data.** This section delves into the analysis of a global food trading dataset initially assembled by De Domenico et al. (2015), accessible at http://www.fao.org. The dataset includes economic networks covering a range of products, where countries are represented as nodes and the edges indicate trade links for particular food products. Following the same preprocessing steps as Jing et al. (2021), we prepared the data to establish a common ground for comparing clustering outcomes. The original directed networks were simplified by omitting their directional features, thereby converting them into undirected networks.
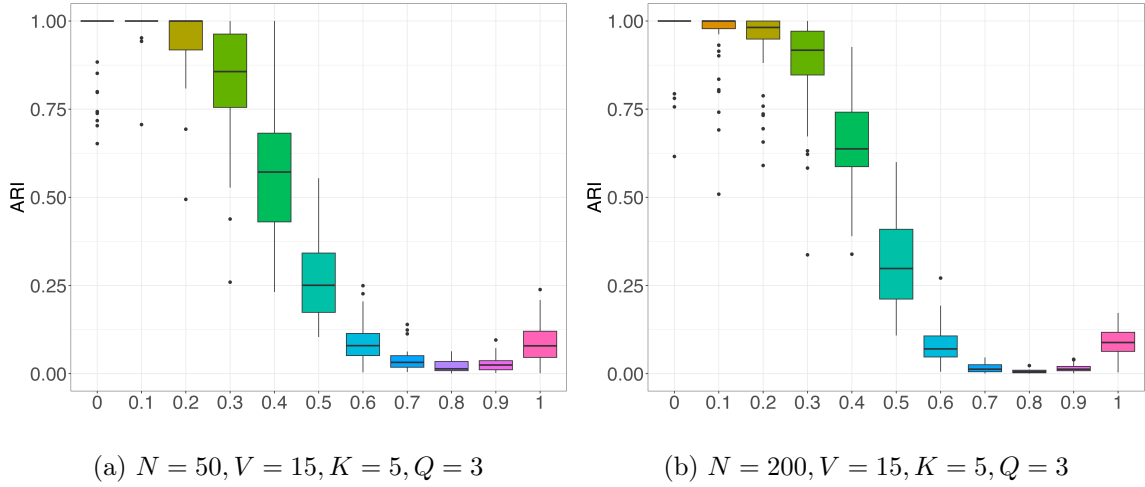
(a) $N = 50, V = 15, K = 5, Q = 3$          (b) $N = 200, V = 15, K = 5, Q = 3$

Figure 10: Performances of *mimi-SBM* on individual clustering through the evolution of label-switching rate.



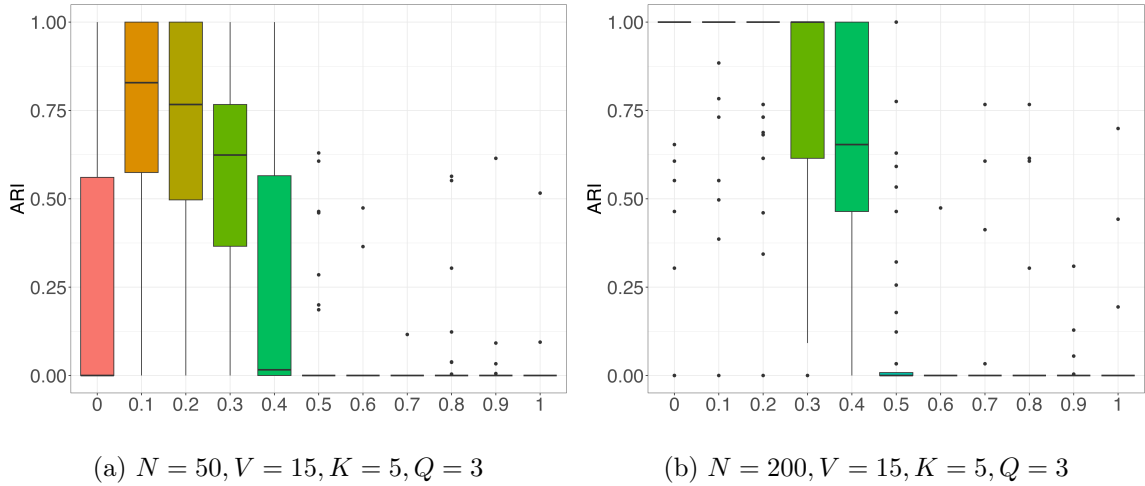(a) $N = 50, V = 15, K = 5, Q = 3$          (b) $N = 200, V = 15, K = 5, Q = 3$

Figure 11: Performances of *mimi-SBM* on view clustering through the evolution of label-switching rate.

Subsequently, to effectively filter out less significant information from the dataset, we eliminate links with a weight of less than 8 and layers containing limited information (less than 150 nodes). Finally, the intersections of the biggest networks of the preselected layers are then extracted. Each layer reflects the international trade interactions involving 30 distinct food products among 99 different countries and regions (nodes).

**TWIST analysis.** In our research, we followed the analytical process described in Jing et al. (2021), to facilitate reliable comparison of results. Consistent with this methodology, we fixed the number of clusters at $K = 4$ for individuals and $Q = 2$ for views.

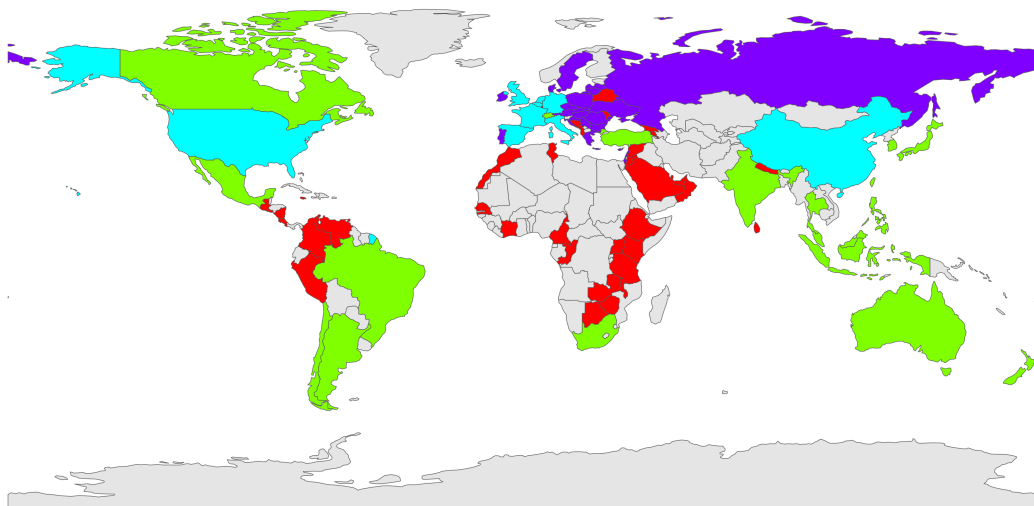In Figure 12, clusters have their own interaction patterns:

Figure 12: World Map of Clusters: Countries are color-coded based on the clusters identified by the model. The cyan cluster (cluster 1) encompasses the West and China; The violet cluster (cluster 2) consists of Russia and some parts of Western Europe; The red cluster (cluster 3) includes countries from Africa and Central America; The green cluster (cluster 4) covers Mexico, Canada, India, Australia, South Africa, Japan, among others; Countries depicted in grey are not included in the database analyzed.

- Cluster 1 serves as a hub due to its centralization of exchanges, exhibiting a high intra-connectivity ($> 90\%$) and substantial inter-connectivity ($> 70\%$), as revealed by the multilayer adjacency probability analysis.

- Cluster 2 displays a robust intra-connectivity, with notable interactions observed with both clusters 1 and 4. Conversely, exchanges with cluster 3 are infrequent for the commodities comprising the database.

- Cluster 3 and Cluster 4 exhibit both intra-cluster and inter-cluster interactions, with a preference for inter-cluster interactions with cluster 1. However, while Cluster 3 predominantly interacts with Cluster 1, Cluster 4 demonstrates partial interaction with Cluster 2.

Exploration of the view components in Table 1 reveals a marked tendency to distinguish between "processed products" and "unprocessed products", although there are some notable exceptions. In addition, it should be noted that Component 1 displays more important connections than Component 2, suggesting that the main flow of transactions is mainly concentrated on products included in Component 1. This observation reinforces Cluster 1's position as a central hub, remaining a predominant actor in the concentration of trade within the various components.

| View component 1 | Beverages_non_alcoholic , Food_prep_nes, Chocolate_products_nes , Crude_materials, Fruit_prepared_nes, Beverages_distilled_alcoholic, Pastry, Sugar_confectionery, Wine |
|---|---|
| View component 2 | Cheese_whole_cow_milk, Cigarettes, Flour_wheat Beer_of_barley, Cereals_breakfast, Coffee_green, Milk_skimmed_dried, Juice_fruit_nes, Maize, Macaroni, Oil_palm, Milk_whole_dried, Oil_essential_nes, Rice_milled, Sugar_refined, Tea Spices_nes, Vegetables_preserved_nes, Water_ice_etc, Vegetables_fresh_nes, Tobacco_unmanufactured |

Table 1: Table of members in view components.

The analysis carried out in this study is reflected in a striking correlation with the steps taken in the precedent analysis. Firstly, we found that the same partitions of individuals were present, with only minor variations in clustering. The links forged within these groups proved to be consistent with market dynamics, highlighting, in particular, the hub role played by cluster 1 in global trade. Furthermore, the overall partitioning of food types persisted, illustrating the persistent distinction between processed and unprocessed products, although a few exceptions were noted, similar to those observed in the previous analysis. In sum, our results largely converge with those of the Jing et al. (2021) study, although a few discrepancies remain, underlining the importance of continuing research in this area to refine our understanding and approach.

**Our optimization.** First, the criterion for choosing the optimal model was employed to guide the selection of hyperparameters. A grid search was conducted over a range of values, spanning from 1 to 20 for the hyperparameter $K$ and from 1 to 10 for $Q$, in concordance with parameters of the first core in Jing et al. (2021) experimentation. The model selection process led to the choice of hyperparameters $K = 20$ and $Q = 1$ as the most suitable configuration. The model found it excessively costly to introduce additional components across the views compared to the information gain achieved, so $Q = 1$ was selected. For individual clustering, $K = 20$ was based on the model's discovery of numerous micro-clusters representing countries based on their interaction habits. This indicates that the model successfully identified fine-grained distinctions among countries, revealing intricate subgroups within the data.

The results show that certain clusters have been substantially preserved, in particular Cluster 1, which remains virtually intact, as does Cluster 4. However, there has been a significant fragmentation of existing clusters, with China in particular remaining isolated. There have also been significant changes in the configuration of clusters, notably the inclusion of Russia along with other South American countries. This change could be explained by the fact that we are now considering only one view component, Russia is closer, in the sense of trading more with, the countries of South America than the rest of the world.
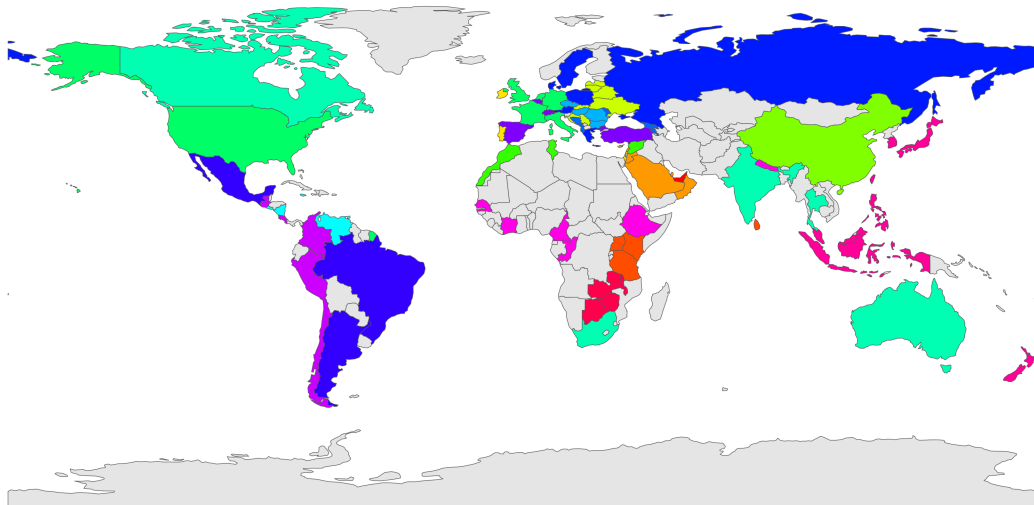
Figure 13: Clustering world map: countries are colored according to the clusters, and parameters defined by the model ($K = 20$).

## Conclusion

This paper proposes a new framework for Mixture of Multilayer SBM *mimi-SBM* that stratifies individuals as well as views.

In order to get a manageable lower bound on the observed log-likelihood, a variational Bayesian approach has been devised. Each model parameter has been estimated using a Variational Bayes EM algorithm. The advantage of such a Bayesian framework consists in allowing the development of an efficient model selection strategy. Moreover, we have provided the proof of model identifiability for the *mimi-SBM* parameters.

In our simulation setting, the *mimi-SBM* related algorithm has been shown to compete with methods based on tensor decomposition, hierarchical model-based SBM, and reference model in consensus clustering in two critical aspects of data analysis: individual clustering and view component identification. Specifically, our algorithm reliably recovered the primary sources of information in the majority of investigated cases. These remarkable performances attest to the efficacy of our approach, underscoring its potential for diverse applications requiring a profound understanding of complex data structures. In real-world application on *Worldwide Food Trading Networks*, when considering the paradigm provided by Jing et al. (2021), we obtain consistent results. However, upon optimizing our model using our metric, a distinctly different clustering emerges. This alternative clustering not only diverges significantly but also reflects much finer nuances in transactional natures.

An interesting follow-up would be to extend this approach to the context of deep learning, specifically in the context of variational auto-encoders using the Bayesian formulation. Additionally, further research is needed to develop theoretical proofs regarding the convergence of parameters for the component-connection probability tensor model.

**Acknowledgments and Disclosure of Funding**

# References

H. Akaike. Information theory and an extension of the maximum likelihood principle. In Selected papers of hirotugu akaike, pages 199–213. Springer, 1998.

H. Attias. A variational baysian framework for graphical models. Advances in neural information processing systems, 12, 1999.

T. Baltrušaitis, C. Ahuja, and L.-P. Morency. Multimodal machine learning: A survey and taxonomy. IEEE Transactions on Pattern Analysis and Machine Intelligence, 41(2): 423–443, 2018.

P. Barbillon, S. Donnet, E. Lazega, and A. Bar-Hen. Stochastic block models for multiplex networks: an application to a multilevel network of researchers. Journal of the Royal Statistical Society Series A: Statistics in Society, 180(1):295–314, 2017.

C. Biernacki, G. Celeux, and G. Govaert. Assessing a mixture model for clustering with the integrated completed likelihood. IEEE transactions on pattern analysis and machine intelligence, 22(7):719–725, 2000.

C. Biernacki, G. Celeux, and G. Govaert. Exact and monte carlo calculations of integrated likelihoods for the latent class model. Journal of Statistical Planning and Inference, 140 (11):2991–3002, 2010.

R. Boutalbi, L. Labiod, and M. Nadif. Tensor latent block model for co-clustering. International Journal of Data Science and Analytics, 10(2):161–175, 2020.

R. Boutalbi, L. Labiod, and M. Nadif. Implicit consensus clustering from multiple graphs. Data Mining and Knowledge Discovery, 35:2313–2340, 2021.

G. Celeux and G. Govaert. A classification em algorithm for clustering and two stochastic versions. Computational statistics & Data analysis, 14(3):315–332, 1992.

A. Celisse, J.-J. Daudin, and L. Pierre. Consistency of maximum-likelihood and variational estimators in the stochastic block model. Electron. J. Statist., 2012.

P.-Y. Chen and A. O. Hero. Multilayer spectral graph clustering via convex layer aggregation: Theory and algorithms. IEEE Transactions on Signal and Information Processing over Networks, 3(3):553–567, 2017.

E. Côme and P. Latouche. Model selection and clustering in stochastic block models based on the exact integrated complete data likelihood. Statistical Modelling, 15(6):564–589, 2015.

A. Cornuéjols, C. Wemmert, P. Gançarski, and Y. Bennani. Collaborative clustering: Why, when, what and how. Information Fusion, 39:81–95, 2018.

J.-J. Daudin, F. Picard, and S. Robin. A mixture model for random graphs. Statistics and computing, 18(2):173–183, 2008.

M. De Domenico, V. Nicosia, A. Arenas, and V. Latora. Structural reducibility of multilayer networks. Nature communications, 6(1):6864, 2015.

Y. Ektefaie, G. Dasoulas, A. Noori, M. Farhat, and M. Zitnik. Multimodal learning with graphs. Nature Machine Intelligence, 5(4):340–350, 2023.

E. Elhamifar and R. Vidal. Sparse subspace clustering: Algorithm, theory, and applications. IEEE transactions on pattern analysis and machine intelligence, 35(11):2765–2781, 2013.

X. Fan, M. Pensky, F. Yu, and T. Zhang. Alma: alternating minimization algorithm for clustering mixture multilayer network. The Journal of Machine Learning Research, 23(1): 14855–14900, 2022.

A. L. Fred and A. K. Jain. Combining multiple clusterings using evidence accumulation. IEEE transactions on pattern analysis and machine intelligence, 27(6):835–850, 2005.

K. Golalipour, E. Akbari, S. S. Hamidi, M. Lee, and R. Enayatifar. From clustering to clustering ensemble selection: A review. Engineering Applications of Artificial Intelligence, 104:104388, 2021.

Q. Han, K. Xu, and E. Airoldi. Consistent estimation of dynamic and multi-layer block models. In International Conference on Machine Learning, pages 1511–1520. PMLR, 2015.

R. Han, Y. Luo, M. Wang, and A. R. Zhang. Exact clustering in tensor block model: Statistical optimality and computational limit. Journal of the Royal Statistical Society Series B: Statistical Methodology, 84(5):1666–1698, 2022.

S. Huang, H. Weng, and Y. Feng. Spectral clustering via adaptive layer aggregation for multi-layer networks. Journal of Computational and Graphical Statistics, pages 1–15, 2022.

L. Hubert and P. Arabie. Comparing partitions. Journal of classification, 2:193–218, 1985.

B.-Y. Jing, T. Li, Z. Lyu, and D. Xia. Community detection on mixture multilayer networks via regularized tensor decomposition. The Annals of Statistics, 49(6):3181–3205, 2021.

P. Latouche, E. Birmele, and C. Ambroise. Variational bayesian inference and complexity control for stochastic block models. Statistical Modelling, 12(1):93–115, 2012.

Y. Li, F. Nie, H. Huang, and J. Huang. Large-scale multi-view spectral clustering via bipartite graph. In Proceedings of the AAAI conference on artificial intelligence, volume 29, 2015.

L. Liu, F. Nie, A. Wiliem, Z. Li, T. Zhang, and B. C. Lovell. Multi-modal joint clustering with application for unsupervised attribute discovery. IEEE Transactions on Image Processing, 27(9):4345–4356, 2018.

P. Mercado, A. Gautier, F. Tudisco, and M. Hein. The power mean laplacian for multilayer graph clustering. In International Conference on Artificial Intelligence and Statistics, pages 1828–1838. PMLR, 2018.

S. Monti, P. Tamayo, J. Mesirov, and T. Golub. Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. Machine Learning, 52:91–118, 2003.

M. Noroozi and M. Pensky. Sparse subspace clustering in diverse multiplex network model. arXiv preprint arXiv:2206.07602, 2022.

S. Paul and Y. Chen. Consistent community detection in multi-relational data through restricted multi-layer stochastic blockmodel. Electronic Journal of Statistics, 10(2):3807 – 3870, 2016.

S. Paul and Y. Chen. Spectral and matrix factorization methods for consistent community detection in multi-layer networks. The Annals of Statistics, 48(1):230 – 250, 2020. doi: 10.1214/18-AOS1800.

M. Pensky and Y. Wang. Clustering of diverse multiplex networks. arXiv preprint arXiv:2110.05308, 2021.

T. Rebafka. Model-based clustering of multiple networks with a hierarchical algorithm, 2023.

G. Schwarz. Estimating the dimension of a model. The annals of statistics, pages 461–464, 1978.

N. Stanley, S. Shai, D. Taylor, and P. J. Mucha. Clustering network layers with the strata multilayer stochastic block model. IEEE transactions on network science and engineering, 3(2):95–105, 2016.

A. Strehl and J. Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. Journal of machine learning research, 3(Dec):583–617, 2002.

U. Von Luxburg. A tutorial on spectral clustering. Statistics and computing, 17:395–416, 2007.

U. Von Luxburg, M. Belkin, and O. Bousquet. Consistency of spectral clustering. The Annals of Statistics, pages 555–586, 2008.

M. Wang and Y. Zeng. Multiway clustering via tensor block models. Advances in neural information processing systems, 32, 2019.

J. Zhao, X. Xie, X. Xu, and S. Sun. Multi-view learning overview: Recent progress and new challenges. Information Fusion, 38:43–54, 2017.

## Appendix A. Identifiability

This appendix is dedicated to the proof of the theorem of Section 3.2 related to the identifiability of the parameters of *mimi-SBM*, recalled below. The proof is very similar to the one of Celisse et al. (2012) and make use of algebraic properties to prove that the parameters depend solely on the marginal distribution of our data.

**Theorem 2** *Let $N \geq \max(2K, 4Q)$ and $V \geq 2K$. Assume that for any $1 \leq k, l \leq K$ and every $1 \leq s \leq Q$, the coordinates of $\boldsymbol{\pi}^T \boldsymbol{\alpha}_{k..} \boldsymbol{\rho}$ are all different, $(\boldsymbol{\pi}^T \boldsymbol{\alpha}_{..s} \boldsymbol{\pi})_{s=1:Q}$ are distinct, and each $(\boldsymbol{\alpha}_{kl.} \boldsymbol{\rho})_{k,l=1:K}$ differs. Then, the mimi-SBM parameter $\boldsymbol{\Theta} = (\boldsymbol{\pi}, \boldsymbol{\rho}, \boldsymbol{\alpha})$ is identifiable.*

### A.1 Assumptions

$\mathcal{A}1$: $(\boldsymbol{\pi}^T \boldsymbol{\alpha}_{k..} \boldsymbol{\rho})_{k=1:K}$ are all different.

$\mathcal{A}2$: $(\boldsymbol{\pi}^T \boldsymbol{\alpha}_{..s} \boldsymbol{\pi})_{s=1:Q}$ are all different.

$\mathcal{A}3$: $N, V \geq 2K$.

$\mathcal{A}4$: $N \geq 4Q$.

$\mathcal{A}5$: $(\boldsymbol{\alpha}_{kl.} \boldsymbol{\rho})_{k,l=1:K}$ are all different.

### A.2 Identifiability of $\boldsymbol{\pi}$

To prove the identifiability of $\boldsymbol{\pi}$, we first need to establish some correspondences. For any $1 \leq k \leq K$, $\forall (i,j,v)$, let $r_k$ be the probability that an edge between $i$ and $j$ in layer $v$ given individual $i$ is in the cluster $k$:

$$
\begin{aligned}
r_k &= \mathbb{P}(A_{ijv} = 1 \mid Z_i = k) \\
&= \sum_l \sum_s \mathbb{P}(A_{ijv} = 1 \mid Z_i = k, Z_j = l, W_v = s)\, \pi_l\, \rho_s \\
&= \sum_l \sum_s \alpha_{kls}\, \pi_l\, \rho_s \\
&= \boldsymbol{\pi}^T \boldsymbol{\alpha}_{k..} \boldsymbol{\rho}\,.
\end{aligned}
$$

**Proposition 3 (Invertibility of R)** *Let $\mathbf{R}$ denote a Vandermonde matrix of size $K \times K$ such as $R_{ik} = (r_k)^{i-1}$, $1 \leq i, k \leq K$. $\mathbf{R}$ is invertible, since the coordinates of $r$ are all different according to Assumption $\mathcal{A}1$.*

Furthermore, for $2 \leq i \leq K$, the joint probability of having $(i - 1)$ edges is given by:

$$
\mathbb{P}(A_{121} = 1, A_{132} = 1, \ldots, A_{1i(i-1)} = 1 \mid Z_1 = k)
$$
$$
= \sum_l \sum_s \mathbb{P}(A_{121} = 1, A_{132} = 1, \ldots, A_{1i(i-1)} = 1 \mid Z_1 = k, Z_2 = l, W_1 = s) \; \pi_l \; \rho_s
$$
$$
= \mathbb{P}(A_{132} = 1, \ldots, A_{1i(1-1)} = 1 \mid Z_1 = k) \times \sum_l \sum_s \mathbb{P}(A_{121} = 1 \mid Z_1 = k, Z_2 = l, W_1 = s) \; \pi_l \; \rho_s
$$
$$
= \mathbb{P}(A_{132} = 1, \ldots, A_{1i(1-1)} = 1 \mid Z_1 = k) \; r_k
$$
$$
= (r_k)^{i-1}
$$
$$
= R_{ik} \, .
$$

Now, we define $u_0 = 1$ and for $1 \leq i \leq 2K - 1$:

$$
u_i = \mathbb{P}(A_{121} = 1, A_{132} = 1, \ldots, A_{1i(1-1)} = 1, A_{1(i+1)i} = 1)
$$
$$
= \sum_k \mathbb{P}(A_{121} = 1, A_{132} = 1, \ldots, A_{1(1+1)i} = 1 \mid Z_1 = k) \; \pi_k
$$
$$
= \sum_k (r_k)^i \; \pi_k \, .
$$

By Assumption $\mathcal{A}3$, $(u_i)_{i=1:(2K-1)}$ are well defined. Hence, $u_0 = 1$ and $(u_i)_{i=1:(2K-1)}$ are known and defined from the marginal $\mathbb{P}_A$. As a consequence, $(u_i)_{i=1:(2K-1)}$ are identifiable.

Also, let $\mathbf{M}$ of size $(K + 1) \times K$ be the matrix given by $M_{ij} = u_{i+j-2}$ for $1 \leq i \leq K + 1$ and $1 \leq j \leq K$, and let $\mathbf{M}_{-i}$ denote the square matrix obtained by removing the row $i$ from $\mathbf{M}$. The coefficients of $\mathbf{M}_{-(K+1)}$, for $1 \leq i, j \leq K$, are:

$$
M_{ij} = \sum_{k=1}^{K} (r_k)^{i-1} \; \pi_k \; (r_k)^{j-1} \, , \text{ and}
$$
$$
\mathbf{M}_{-(K+1)} = \mathbf{R} \operatorname{Diag}(\boldsymbol{\pi}) \mathbf{R}^T \, . \tag{11}
$$

**Proposition 4 (Relations between R, M and $\boldsymbol{\pi}$)** *From Proposition 3 and Equation (11), we can define*

$$
\mathbf{M}_{-(K+1)} = \mathbf{R} \; \operatorname{Diag}(\boldsymbol{\pi}) \; \mathbf{R}^T \, .
$$

The correspondence of the different terms being established, we now need to prove the identifiability of $\boldsymbol{\pi}$, which means showing that $\mathbf{M}_{-(K+1)}$ and $\mathbf{R}$ are identifiable.

First, for the identifiability of $r_k$, with $\delta_k = \operatorname{Det}(\mathbf{M}_{-k})$, we define a polynomial function $B$ such as:

$$
B(x) = \sum_{k=0}^{K} (-1)^{K+k} \; \delta_{k+1} \; x^k \, .
$$

This polynomial function has two important properties.

**Proposition 5** *Let* $\deg(B)$ *denote the degree of B. We have* $\deg(B) = K$.

**Proof** Let $\delta_{K+1} = \text{Det}(\mathbf{M}_{-(K+1)})$, with $M_{-(K+1)} = \mathbf{R} \ \text{Diag}(\boldsymbol{\pi}) \ \mathbf{R}^T$ as stated in Proposition 4, and $\mathbf{R}$ being invertible as stated in Proposition 3. In consequence, $\mathbf{M}_{-(K+1)}$ is the product of invertible matrices, $\delta_{K+1} = \text{Det}(\mathbf{M}_{-(K+1)}) \neq 0$ and, moreover, $\deg(B) = K$. ∎

**Proposition 6** *For* $1 \leq k \leq K$, $B(r_k) = 0$.

**Proof** Let $\mathbf{N}_k$ of size $(K+1) \times (K+1)$ be the concatenation in columns of the matrix $\mathbf{M}$ with the vector $V_k = [1, r_k, r_k^2, \ldots, r_k^K]^T$.

Now let's calculate the determinant of $\mathbf{N}_k$ developed by the last column:

$$
\begin{aligned}
\det(\mathbf{N}_k) &= \sum_{l=0}^{K} (-1)^{K+1+l+1} (r_k)^l \ \det(\mathbf{M}_{l+1}) \\
&= \sum_{l=0}^{K} (-1)^{K+l} \delta_{l+1} (r_k)^l \\
&= B(r_k).
\end{aligned}
$$

In addition, the $j$th column of the $\mathbf{M}$ matrix can be written as $M_{\cdot j} = \sum_{k=1}^{K} r_k^{j-1} \pi_k V_k$. Therefore, $\text{rank}(\mathbf{N}_k) < K+1$ and $\det(\mathbf{N}_k) = 0$ for $1 \leq k \leq K$. In consequence, $B(r_k) = 0$ for $1 \leq k \leq K$. ∎

With $(r_k)_{k=1:K}$ being the roots of $B$ (proposition 6), they are functions of $(\delta_k)_{k=1:K+1}$ which are themselves derived from $\mathbb{P}_A$. Also, $(r_k)_{k=1:K}$ can be expressed in a unique way (up to label switching) from $\mathbb{P}_A$, thus $(r_k)_{k=1:K}$ are identifiable. In consequence, $\mathbf{R}$ is also identifiable by definition. Finally, since $\mathbf{M}_{-(K+1)}$ and $\mathbf{R}$ are identifiable and invertible, $\text{Diag}(\boldsymbol{\pi}) = \mathbf{R}^{-1}\mathbf{M}_{-(K+1)}(\mathbf{R}^T)^{-1}$. In conclusion, $\boldsymbol{\pi}$ is identifiable.

### A.3 Identifiability of $\rho$

Identifiability of $\boldsymbol{\rho}$ is similar to $\boldsymbol{\pi}$, the main difference lies in the assumptions made and the quantities defined.

For any $1 \leq s \leq Q$, $\forall(i, j, v)$, let $t_s$ be the probability of an edge between $i$ and $j$ in layer $v$ given view $v$ is in the component $s$:

$$
\begin{aligned}
t_s &= \mathbb{P}(A_{ijv} = 1 \mid W_v = s) \\
&= \sum_{l} \sum_{k} \mathbb{P}(A_{ijv} = 1 \mid Z_i = k, Z_j = l, W_v = s) \ \pi_l \ \pi_k \\
&= \sum_{l} \sum_{k} \alpha_{kls} \ \pi_l \ \pi_k \\
&= \boldsymbol{\pi}^T \boldsymbol{\alpha}_{\cdot \cdot s} \boldsymbol{\pi} .
\end{aligned}
$$

**Proposition 7 (Invertibility of T)** *Let* $\mathbf{T}$ *denote a Vandermonde matrix of size* $Q \times Q$ *such as* $T_{is} = (t_s)^{i-1}$, $1 \leq i, s \leq Q$. $\mathbf{T}$ *is invertible, since the coordinates of* $(t_s)$ *are all different according to Assumption* $\mathcal{A}2$.

Let's define the joint probability of $i-1$ edges given the latent component of the view 1:

$$\mathbb{P}(A_{121} = 1, A_{341} = 1, \ldots, A_{2i-1\,2i\,1} = 1 \mid W_1 = s)$$

$$= \sum_l \sum_k \mathbb{P}(A_{121} = 1, A_{341} = 1, \ldots, A_{2i-1\,2i\,1} \mid Z_1 = k, Z_2 = l, W_1 = s)\ \pi_l\ \pi_k$$

$$= \mathbb{P}(A_{341} = 1, \ldots, A_{2i-1\,2i\,1} = 1 \mid W_1 = s) \times \sum_l \sum_k \mathbb{P}(A_{121} = 1 \mid Z_1 = k, Z_2 = l, W_1 = s)\ \pi_l\ \pi_k$$

$$= \mathbb{P}(A_{341} = 1, \ldots, A_{2i-1\,2i\,1} = 1 \mid W_1 = s) \times\ t_s$$

$$= (t_s)^{i-1}.$$

Now, we define $v_0 = 1$ and for $1 \leq i \leq 2Q - 1$:

$$v_i = \mathbb{P}(A_{121} = 1, A_{341} = 1, \ldots, A_{2i\,2i+1\,1} = 1)$$

$$= \sum_s \mathbb{P}(A_{121} = 1, A_{341} = 1, \ldots, A_{2i\,2i+1\,1} = 1 \mid W_1 = s)\ \rho_s$$

$$= \sum_s (t_s)^i\ \rho_s.$$

By Assumption $\mathcal{A}4$, $(v_i)_{i=1:(2Q-1)}$ are well defined. Hence, $v_0 = 1$ and $(v_i)_{i=1:(2Q-1)}$ are known and defined from the marginal $\mathbb{P}_A$. As a consequence, $(v_i)_{i=1:(2Q-1)}$ are identifiable.

Also, let $\tilde{\mathbf{M}}$ be the matrix of size $(Q+1) \times Q$ given by $\tilde{M}_{ij} = v_{i+j-2}$ for $1 \leq i \leq Q+1$ and $1 \leq j \leq Q$, and let $\tilde{\mathbf{M}}_{-i}$ denote the square matrix obtained by removing the row $i$ from $\tilde{\mathbf{M}}$. The coefficients of $\tilde{\mathbf{M}}_{-(K+1)}$, for $1 \leq i, j \leq Q$, are:

$$\tilde{M}_{ij} = \sum_{s=1}^{Q} (t_s)^{i-1}\ \rho_s\ (r_s)^{j-1}, \ \text{and}$$

$$\tilde{\mathbf{M}}_{-(Q+1)} = \mathbf{T}\ \mathrm{Diag}(\boldsymbol{\rho})\ \mathbf{T}^T. \tag{12}$$

**Proposition 8 (Relations between T, $\tilde{\mathbf{M}}$ and $\boldsymbol{\rho}$)** *From Proposition 7 and Equation* (12), *we can define*

$$\tilde{\mathbf{M}}_{-(Q+1)} = \mathbf{T}\ \mathrm{Diag}(\boldsymbol{\rho})\ \mathbf{T}^T.$$

The correspondence of the different terms being established, we now need to prove the identifiability of $\boldsymbol{\rho}$, which means showing that $\tilde{\mathbf{M}}_{-(Q+1)}$ and $\mathbf{T}$ are identifiable.

As for the previous proof regarding the identifiability of $t_s$, with $\delta_s = \mathrm{Det}(\tilde{\mathbf{M}}_{-s})$, we define a polynomial function $\tilde{B}$ such as:

$$\tilde{B}(x) = \sum_{s=0}^{Q} (-1)^{Q+s}\ \delta_{s+1}\ x^s$$

This polynomial function has again two important properties summarized in the following proposition.

**Proposition 9** *Let* $\deg(\tilde{B})$ *denote the degree of* $\tilde{B}$. *We have* $\deg(\tilde{B}) = Q$ *and* $\tilde{B}(t_s) = 0$, *for* $1 \leq s \leq Q$.

**Proof** The proof follow the same lines as those of Proposition 5 and Proposition 6. ∎

With $(t_s)_{s=1:Q}$ being the roots of $\tilde{B}$ (proposition 9), they are functions of $(\delta_s)_{s=1:Q+1}$ which are themselves derived from $\mathbb{P}_A$. Also, $(t_s)_{s=1:Q}$ can be expressed in a unique way (up to label switching) from $\mathbb{P}_A$, thus $(t_s)_{s=1:Q}$ are identifiable. In consequence, $\mathbf{T}$ is also identifiable by definition. Finally, since $\tilde{\mathbf{M}}_{-(Q+1)}$ and $\mathbf{T}$ are identifiable and invertible, $\text{Diag}(\boldsymbol{\rho}) = \mathbf{T}^{-1}\tilde{\mathbf{M}}_{-(Q+1)}(\mathbf{T}^T)^{-1}$. In conclusion, $\boldsymbol{\rho}$ is identifiable.

### A.4 Identifiability of $\boldsymbol{\alpha}$

To establish the identifiability of $\boldsymbol{\alpha}$, the initial proof relies on matrix inversion. However, within our framework, tensor inversion is not as straightforward as uniqueness may not be inherently guaranteed. To overcome this issue, we will reparametrize our problem to revert to a matrix-based formulation. To do this, we shift from utilizing the reference frame of nodes (individuals) to that of edges (connections).

First, the tensor $\mathbf{A}$ is transformed into a matrix $\tilde{\mathbf{A}}$ of size $\tilde{N} \times Q$, with $\tilde{N} = N(N-1)/2$ in an undirected framework. Each column corresponds to a vectorization of the upper triangular matrix of each layer of $\mathbf{A}$. Thus, the edge $\mathbf{A}_{ijv}$ will be described by $\tilde{\mathbf{A}}_{\tilde{i}v}$, with $\tilde{i}$ being the index corresponding to the edge $(i,j)$ between nodes $i$ and $j$.

Then, we can map the clustering of observations into a clustering of edges, which results in a matrix $\tilde{\mathbf{Z}}$ of size $\tilde{N} \times \tilde{K}$, with $\tilde{K} = K(K+1)/2$. Each row of the matrix corresponds to the clustering of the pair of nodes making up the edges $1 \leq \tilde{i} \leq \tilde{N}$.

Also, we denote $\tilde{\boldsymbol{\pi}}$ the proportion vector of pairs such that $\tilde{\boldsymbol{\pi}}_{\tilde{k}} = \boldsymbol{\pi}_k\boldsymbol{\pi}_l$, for $1 \leq \tilde{k} \leq \tilde{K}$, where $1 \leq k, l \leq K$ are the initial clusters corresponding to the index of the $\tilde{k}$ in the reparametrization.
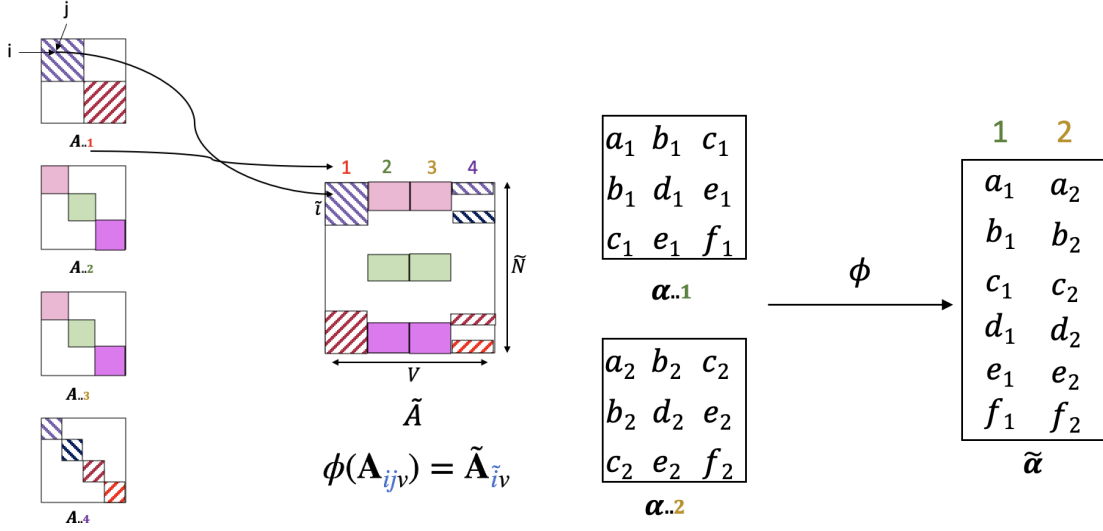
Finally, $\tilde{\boldsymbol{\alpha}}$ is a $\tilde{K} \times Q$ matrix whose rows represents the clusters related to the pairs of nodes while the columns are the components. The terms of $\tilde{\boldsymbol{\alpha}}$ represent the probabilities of connection between these clusters and components.

Now, let's define a function $\phi$ such as:

$$\phi(\mathbf{A}, \mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\alpha}) = (\tilde{\mathbf{A}}, \tilde{\mathbf{Z}}, \tilde{\boldsymbol{\pi}}, \tilde{\boldsymbol{\alpha}}).$$

The function is bijective for $\mathbf{A}$ and $\boldsymbol{\alpha}$ and injective for $\mathbf{Z}$ and $\boldsymbol{\pi}$ ; The bijective relationship involving the parameter $\boldsymbol{\alpha}$ and $\tilde{\boldsymbol{\alpha}}$ enables the establishment of identifiability. The aim is therefore to show the identifiability of $\tilde{\boldsymbol{\alpha}}$.

**Remark 10** *These transformations map our problem into a LBM framework (see Figure 14). Hence, the identifiability of* $\tilde{\boldsymbol{\alpha}}$ *will be developed accordingly.*

Figure 14: Illustration of transformation on $\mathbf{A}$ and $\boldsymbol{\alpha}$.

The proof is identical to the ones of Section A.2 For any $1 \leq \tilde{k} \leq \tilde{K}$ and $\forall i, j$, let's define:

$$
\begin{aligned}
\tilde{r}_{\tilde{k}} &= \mathbb{P}(\tilde{A}_{ij} = 1 \mid \tilde{Z}_i = \tilde{k}) \\
&= \sum_s \mathbb{P}(\tilde{A}_{ij} = 1 \mid \tilde{Z}_i = \tilde{k}, W_v = s) \, \rho_s \\
&= \sum_s \tilde{\alpha}_{\tilde{k}s} \rho_s \\
&= (\tilde{\boldsymbol{\alpha}} \boldsymbol{\rho})_{\tilde{k}} \, .
\end{aligned}
$$

**Proposition 11 (Invertibility of $\tilde{\mathbf{R}}$)** *Let $\tilde{\mathbf{R}}$ denote a Vandermonde matrix of size $\tilde{K} \times \tilde{K}$ such as $\tilde{R}_{i\tilde{k}} = (\tilde{r}_{\tilde{k}})^{i-1}$, for $1 \leq i \leq \tilde{K}$ and $1 \leq \tilde{k} \leq \tilde{K}$. $\tilde{\mathbf{R}}$ is invertible, since the coordinates of $r$ are all different according to Assumption $\mathcal{A}5$.*

The rest of the proof is identical to the one of Section A.2 so that it can be show that $\tilde{\mathbf{R}}$ is identifiable and so $\tilde{\boldsymbol{\pi}}_{\tilde{k}}$.

We now focus on the the identifiability of $\tilde{\boldsymbol{\alpha}}$. Let $\mathbf{U}$ be a matrix of size $\tilde{K} \times Q$ such that the $(i, j)$ entry of is the joint probability of having $i$ connections in the first row and $j - 1$ connections in the first column:

$$
\begin{aligned}
\mathbf{U}_{ij} &= \mathbb{P}(\tilde{\mathbf{A}}_{11} = 1, \tilde{\mathbf{A}}_{12} = 1, \ldots, \tilde{\mathbf{A}}_{1i} = 1, \tilde{\mathbf{A}}_{21} = 1, \ldots, \tilde{\mathbf{A}}_{j1} = 1) \\
&= \sum_{\tilde{k}} \sum_q \tilde{\pi}_{\tilde{k}} \, \rho_s \, \mathbb{P}(\tilde{\mathbf{A}}_{11} = 1, \tilde{\mathbf{A}}_{12} = 1, \ldots, \tilde{\mathbf{A}}_{1i} = 1, \tilde{\mathbf{A}}_{21} = 1, \ldots, \tilde{\mathbf{A}}_{j1} = 1 \mid \tilde{\mathbf{Z}}_1 = \tilde{k}, \mathbf{W}_1 = s) \\
&= \sum_{\tilde{k}} \sum_q \tilde{\pi}_{\tilde{k}} \, \rho_s \, \tilde{\alpha}_{\tilde{k}s} \, \mathbb{P}(\tilde{\mathbf{A}}_{12} = 1, \ldots, \tilde{\mathbf{A}}_{1i} = 1, \tilde{\mathbf{A}}_{21} = 1, \ldots, \tilde{\mathbf{A}}_{j1} = 1 \mid \tilde{\mathbf{Z}}_1 = \tilde{k}, \mathbf{W}_1 = s) \\
&= \sum_{\tilde{k}} \sum_q \tilde{\pi}_{\tilde{k}} \, \rho_s \, \tilde{\alpha}_{\tilde{k}s} \, \tilde{r}_{\tilde{k}}^{i-1} t_s^{j-1}. \quad (13)
\end{aligned}
$$

34

**Proposition 12 (Relations between $\tilde{\mathbf{R}}$, T, U, $\tilde{\boldsymbol{\alpha}}$, $\tilde{\boldsymbol{\pi}}$ and $\boldsymbol{\rho}$)** *From Proposition 11 and Equation* (13), *we can define* $\mathbf{U} = \tilde{\mathbf{R}} \operatorname{Diag}(\tilde{\boldsymbol{\pi}}) \tilde{\boldsymbol{\alpha}} \operatorname{Diag}(\boldsymbol{\rho}) \mathbf{T}^T$, *with* $\mathbf{U}$, $\tilde{\mathbf{R}}$, $\operatorname{Diag}(\tilde{\boldsymbol{\pi}})$, $\operatorname{Diag}(\boldsymbol{\rho})$ *and* $\mathbf{T}$ *being invertible. Therefore,*

$$\tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{R}}^{-1} \operatorname{Diag}(\tilde{\boldsymbol{\pi}})^{-1} \mathbf{U} \operatorname{Diag}(\rho)^{-1} (\mathbf{T}^T)^{-1}.$$

In addition to Proposition 12, $\mathbf{U}$ being defined from $\mathbb{P}_{\tilde{\mathbf{A}}}$, all its coefficients are identifiable. As a consequence, $\tilde{\boldsymbol{\alpha}}$ is identifiable. In conclusion, $\boldsymbol{\alpha} = \phi^{-1}(\tilde{\boldsymbol{\alpha}})$ is identifiable.

## Appendix B. Details of VBEM algorithm

### B.1 Variational parameters of clustering $\tau_{ik}$

The optimal approximation for $q(\mathbf{Z}_i)$ is

$$q(\mathbf{Z}_i) = \mathcal{M}(\mathbf{Z}_i; (\tau_{i1}, \dots, \tau_{iK})),$$

where $\tau_{ik}$ is the probability of node $i$ to belong to class $k$. It satisfies the relation

$$\tau_{ik} \propto e^{\psi(\beta_k) - \psi(\sum_{k'} \beta_{k'})} \prod_{j \neq i}^{N} \prod_{l=1}^{K} \prod_{v=1}^{V} \prod_{s=1}^{Q} e^{\tau_{jl} \nu_{vs} \left[ A_{ijv} \left( \psi(\eta_{kls}) - \psi(\xi_{kls}) \right) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right]},$$

where $\psi$ is digamma function. Distribution $q(\mathbf{Z})$ is optimized with a fixed point algorithm.
**Proof** According to the model, the optimal distribution $q(\mathbf{Z}_i)$ is given by

$$\log q(\mathbf{Z}_i) = \mathbb{E}_{\mathbf{Z}^{\setminus i}, \alpha, \pi, W, \rho} [\log \mathbb{P}(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \alpha, \pi, \rho)]$$

$$\propto \mathbb{E}_{\mathbf{Z}^{\setminus i}, \alpha, \mathbf{W}} [\log \mathbb{P}(\mathbf{A}|\mathbf{Z}, \mathbf{W}, \alpha)] + \mathbb{E}_{\mathbf{Z}^{\setminus i}, \pi} [\log \mathbb{P}(\mathbf{Z}|\pi)]$$

$$\propto \mathbb{E}_{\mathbf{Z}^{\setminus i}, \alpha, \mathbf{W}} \left[ \sum_{i'=1, j>i'}^{N} \sum_{k,l=1}^{K} \sum_{v=1}^{V} \sum_{s=1}^{Q} \mathbb{1}_{\mathbf{Z}_{i'}, \mathbf{Z}_j, \mathbf{W}_v} \left( \log \mathbb{P}(A_{i'jv}|Z_{i'} = k, Z_j = l, W_v = s, \alpha) \right) \right]$$

$$+ \mathbb{E}_{\mathbf{Z}^{\setminus i}, \pi} \left[ \sum_{i'=1}^{N} \sum_{k}^{K} \log \mathbb{P}(Z_{i'} = k|\pi) \right]$$

$$\propto \sum_{k} \mathbb{1}_{Z_i=k} \left\{ \mathbb{E}_{\pi}[\log(\pi_k)] + \sum_{j \neq i}^{N} \sum_{l=1}^{K} \sum_{v=1}^{V} \sum_{s=1}^{Q} \tau_{jl} \nu_{vs} \mathbb{E}_{\alpha} \left[ A_{ijv} \log(\alpha_{kls}) + (1 - A_{ijv}) \log(1 - \alpha_{kls}) \right] \right\}$$

Remember that :

- $\pi \sim Dir(\pi; \beta)$, so $\pi_k \sim Beta(\pi_k; \beta_k, \sum_{k'} \beta_{k'} - \beta_k)$ ;
- $\mathbb{E}_{\pi}[\log(\pi_k)] = \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'})$;
- $q(\alpha_{kls}) = Beta(\alpha_{kls}; \eta_{kls}, \xi_{kls})$ ;
- $\mathbb{E}_{\alpha_{kls}}[\log(\alpha_{kls})] = \psi(\eta_{kls}) - \psi(\xi_{kls} + \eta_{kls})$;
- $\mathbb{E}_{\alpha_{kls}}[\log(1 - \alpha_{kls})] = \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})$.

In consequence,

$$\log q(Z_i) \propto \sum_{k} \mathbb{1}_{Z_i=k} \left\{ \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) + \sum_{j \neq i}^{N} \sum_{l=1}^{K} \sum_{v=1}^{V} \sum_{s=1}^{Q} \tau_{jl} \nu_{vs} \left[ A_{ijv} \left( (\psi(\eta_{kls}) - \psi(\xi_{kls} + \eta_{kls})) \right. \right. \right.$$

$$- (\psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})) \Big) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \Big] \Big\}$$

$$= \sum_{k} \mathbb{1}_{Z_i=k} \left\{ \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) + \sum_{j \neq i}^{N} \sum_{l=1}^{K} \sum_{v=1}^{V} \sum_{s=1}^{Q} \tau_{jl} \nu_{vs} \left[ A_{ijv} \left( \psi(\eta_{kls}) - \psi(\xi_{kls}) \right) \right. \right.$$

$$+ \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \Big] \Big\}.$$

We can therefore deduce that, by applying the exponential :

$$q(Z_i = k) \propto e^{\psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) + \sum_{j \neq i}^{N} \sum_{l=1}^{K} \sum_{v=1}^{V} \sum_{s=1}^{Q} \tau_{jl} \nu_{vs} \left[ A_{ijv} \left( \psi(\eta_{kls}) - \psi(\xi_{kls}) \right) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right]}$$

$$= e^{\psi(\beta_k) - \psi(\sum_{k'} \beta_{k'})} \prod_{j \neq i}^{N} \prod_{l=1}^{K} \prod_{v=1}^{V} \prod_{s=1}^{Q} e^{\tau_{jl} \nu_{vs} \left[ A_{ijv} \left( \psi(\eta_{kls}) - \psi(\xi_{kls}) \right) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right]}$$ .

Therefore,

$$\tau_{ik} \propto e^{\psi(\beta_k) - \psi(\sum_{k'} \beta_{k'})} \prod_{j \neq i}^{N} \prod_{l=1}^{K} \prod_{v=1}^{V} \prod_{s=1}^{Q} e^{\tau_{jl} \nu_{vs} \left[ A_{ijv} \left( \psi(\eta_{kls}) - \psi(\xi_{kls}) \right) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right]}.$$

So,

$$q(\mathbf{Z}_i) = \mathcal{M}(\mathbf{Z}_i; (\tau_{i1}, \ldots, \tau_{iK})).$$

■

## B.2 Variational parameters of component membership $\nu_{vs}$

The optimal approximation for $q(\mathbf{W}_v)$ is

$$q(W_v) = \mathcal{M}(W_v; (\nu_{v1}, \ldots, \nu_{vQ})),$$

with

$$\nu_{vs} \propto e^{\psi(\theta_s) - \psi(\sum_{s'} \theta_{s'})} \prod_{i \neq j}^{N} \prod_{k \neq l}^{K} e^{\tau_{ik} \tau_{jl} \left[ A_{ijv} \left( \psi(\eta_{kls}) - \psi(\xi_{kls}) \right) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right]}$$

$$\prod_{k}^{K} \prod_{i<j}^{N} e^{\tau_{ik} \tau_{jk} \left[ A_{ijv} \left( \psi(\eta_{kks}) - \psi(\xi_{kks}) \right) + \psi(\xi_{kks}) - \psi(\eta_{kks} + \xi_{kks}) \right]}.$$

$\nu_{vs}$ is the probability of layer $v$ to belong to component $s$.

**Proof** As previously mentioned, in accordance with the principles of variational Bayes, the optimal probability distribution can be expressed as follows:

$$
\begin{aligned}
\log q(\mathbf{W}_v) &= \mathbb{E}_{\mathbf{W}^{\backslash v}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \mathbf{Z}, \boldsymbol{\rho}}[\log \mathbb{P}(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})] \\
&\propto \mathbb{E}_{\mathbf{W}^{\backslash v}, \boldsymbol{\alpha}, \mathbf{Z}}[\log \mathbb{P}(\mathbf{A}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha})] + \mathbb{E}_{\mathbf{W}^{\backslash v}, \boldsymbol{\rho}}[\ln \mathbb{P}(\mathbf{W}|\rho)] \\
&\propto \mathbb{E}_{W^{\backslash v}, \boldsymbol{\alpha}, \mathbf{Z}}\Big[ \sum_{i=1, j>i}^{N} \sum_{k,l=1}^{K} \sum_{v=1}^{V} \sum_{s=1}^{Q} \mathbb{1}_{\mathbf{Z}_i, \mathbf{Z}_j, W_v}\Big( \log \mathbb{P}(A_{ijv}|Z_i = k, Z_j = l, W_v = s, \boldsymbol{\alpha})\Big)\Big] \\
&+ \mathbb{E}_{\mathbf{W}^{\backslash v}, \boldsymbol{\rho}}\Big[\sum_{v=1}^{V} \sum_{s}^{Q} \log \mathbb{P}(W_v = s|\boldsymbol{\rho})\Big] \\
&\propto \sum_q \mathbb{1}_{W_v=s}\Big\{ \mathbb{E}_{\rho}[\log(\rho_s)] + \sum_{k \neq l}^{K} \sum_{i=1, j \neq i}^{N} \tau_{ik}\, \tau_{jl}\, \mathbb{E}_{\boldsymbol{\alpha}}\Big[A_{ijv}\log(\alpha_{kls}) + (1 - A_{ijv})\log(1 - \alpha_{kls})\Big] \\
&+ \sum_k^{K} \sum_{i=1, i<j}^{N} \tau_{ik}\, \tau_{jk}\, \mathbb{E}_{\boldsymbol{\alpha}}\Big[A_{ijv}\log(\alpha_{kks}) + (1 - A_{ijv})\log(1 - \alpha_{kks})\Big]\Big\}.
\end{aligned}
$$

Reminder :

- $\rho \sim Dir(\pi; \theta)$, so $\rho_s \sim Beta(\rho_s; \theta_s, \sum_{s'} \theta_{s'} - \theta_s)$;

- $\mathbb{E}_{\rho}[\log(\rho_s)] = \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'})$ .

Hence,

$$
\begin{aligned}
\log q(W_v) &\propto \sum_q \mathbb{1}_{W_v=q}\Big\{ \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) + \sum_{i=1, j>i}^{N} \sum_{k,l=1}^{K} \tau_{ik}\, \tau_{jl}\Big[A_{ijv}\Big((\psi(\eta_{kls}) - \psi(\xi_{kls} + \eta_{kls})) \\
&- (\psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}))\Big) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})\Big]\Big\} \\
&= \sum_q \mathbb{1}_{W_v=q}\Big\{ \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) + \sum_{i=1, j>i}^{N} \sum_{k,l=1}^{K} \tau_{ik}\, \tau_{jl}\Big[A_{ijv}\Big(\psi(\eta_{kls}) - \psi(\xi_{kls})\Big) \\
&+ \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})\Big]\Big\}.
\end{aligned}
$$

Consequently,

$$
\begin{aligned}
q(W_v = q) &\propto e^{\psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) + \sum_{i=1, j>i}^{N} \sum_{k,l=1}^{K} \tau_{ik}\, \tau_{jl}\Big[A_{ijv}\Big(\psi(\eta_{kls}) - \psi(\xi_{kls})\Big) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})\Big]} \\
&= e^{\psi(\theta_s) - \psi(\sum_{s'} \theta_{s'})} \prod_{k \neq l}^{K} \prod_{i=1, j \neq i}^{N} e^{\tau_{ik}\, \tau_{jl}\Big[A_{ijv}\Big(\psi(\eta_{kls}) - \psi(\xi_{kls})\Big) + \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})\Big]} \\
&\prod_{k}^{K} \prod_{i<j}^{N} e^{\tau_{ik}\, \tau_{jk}\Big[A_{ijv}\Big(\psi(\eta_{kks}) - \psi(\xi_{kks})\Big) + \psi(\xi_{kks}) - \psi(\eta_{kks} + \xi_{kks})\Big]}.
\end{aligned}
$$

So,

$$\nu_{vs} \propto e^{\psi(\theta_s)-\psi(\sum_{s'}\theta_{s'})} \prod_{\substack{i \neq j \\ k \neq l}}^{N} \prod_{k \neq l}^{K} e^{\tau_{ik}\,\tau_{jl}\left[A_{ijv}\left(\psi(\eta_{kls})-\psi(\xi_{kls})\right)+\psi(\xi_{kls})-\psi(\eta_{kls}+\xi_{kls})\right]}$$

$$\prod_{k}^{K} \prod_{i<j}^{N} e^{\tau_{ik}\,\tau_{jk}\left[A_{ijv}\left(\psi(\eta_{kks})-\psi(\xi_{kks})\right)+\psi(\xi_{kks})-\psi(\eta_{kks}+\xi_{kks})\right]}$$ ,

and

$$q(W_v) = \mathcal{M}(W_v; (\nu_{v1}, \ldots, \nu_{vQ})).$$

∎

## B.3 Optimization of $q(\boldsymbol{\pi})$ $(\beta_k)$

Due to the selection of prior distributions, the distribution $q(\boldsymbol{\pi})$ remains within the same family of distributions as the prior distribution $\mathbb{P}(\boldsymbol{\pi})$.

$$q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}; \boldsymbol{\beta}),$$

with

$$\beta_k = \beta_k^0 + \sum_i^N \tau_{ik}.$$

**Proof** The optimal probability distribution can be formulated in the following manner:

$$\begin{aligned} \log q(\boldsymbol{\pi}) &\propto \mathbb{E}_{\mathbf{W},\boldsymbol{\alpha},\mathbf{Z},\boldsymbol{\rho}}[\log \mathbb{P}(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})] \\ &\propto \mathbb{E}_{\mathbf{Z}}[\log \mathbb{P}(\mathbf{Z} \mid \boldsymbol{\pi})] + \log p(\boldsymbol{\pi}) \\ &\propto \sum_i^N \sum_k^K \tau_{ik} \log \pi_k + \sum_{k=1}^K \left(\beta_k^0 - 1\right) \log \pi_k \ . \\ &\propto \sum_k^K \left(\beta_k^0 + (\sum_i^N \tau_{ik}) - 1\right) \log \pi_k \end{aligned}$$

After exponentiation and normalization, we obtain:

$$q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}; \boldsymbol{\beta}),$$

with

$$\beta_k = \beta_k^0 + \sum_i^N \tau_{ik}.$$

∎

**B.4 Optimization of $q(\boldsymbol{\rho})$ $(\theta_s)$**

As previously mentioned, the selection of prior distributions enables us to remain within the same family of distributions.

$$q(\boldsymbol{\rho}) = \text{Dir}(\boldsymbol{\rho}; \boldsymbol{\theta}),$$

with

$$\theta_s = \theta_s^0 + \sum_{v=1}^{V} \nu_{vs}.$$

**Proof** According to variational Bayes, the optimal probability distribution can be expressed as follows:

$$\log q(\boldsymbol{\rho}) \propto \mathbb{E}_{\mathbf{W}, \boldsymbol{\alpha}, Z}[\log \mathbb{P}(\mathbf{A}, \mathbf{Z}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi}, \boldsymbol{\rho})]$$
$$\propto \mathbb{E}_{\mathbf{W}}[\log p(\mathbf{W} \mid \boldsymbol{\rho})] + \log \mathbb{P}(\boldsymbol{\rho})$$
$$\propto \sum_{v}^{V} \sum_{s}^{Q} \nu_{vs} \log \rho_s + \sum_{q=1}^{Q} \left(\theta_s^0 - 1\right) \log \rho_s.$$
$$\propto \sum_{s}^{Q} \left(\theta_s^0 + \left(\sum_{v}^{V} \nu_{vs}\right) - 1\right) \log \rho_s$$

After exponentiation and normalization, we have

$$q(\boldsymbol{\rho}) = \text{Dir}(\boldsymbol{\rho}; \boldsymbol{\theta}),$$

with

$$\theta_s = \theta_s^0 + \sum_{v=1}^{V} \nu_{vs}.$$

∎

**B.5 Optimization of $q(\boldsymbol{\alpha})$ $(\eta_{kls}$ and $\xi_{kls})$**

Once again, the distribution form of the prior distribution $\mathbb{P}(\boldsymbol{\alpha})$ is preserved through the variational optimization process.

$$q(\alpha_{kls}) = \text{Beta}(\alpha_{kls}; \eta_{kls}, \xi_{kls}).$$

When $k \neq l$, parameters $\eta_{kls}$ and $\xi_{kls}$ are given by:

$$\eta_{kls} = \eta_{kls}^0 + \sum_{i \neq j}^{N} \sum_{v}^{V} \tau_{ik} \tau_{jl} \nu_{vs} A_{ijv}$$
$$\xi_{kls} = \xi_{kls}^0 + \sum_{i \neq j}^{N} \sum_{v}^{V} \tau_{ik} \tau_{jl} \nu_{vs} \left(1 - A_{ijv}\right)$$

.

Otherwise, when $k$ equals $l$, the parameters $\eta_{kks}$ and $\xi_{kks}$ are determined by:

$$\eta_{kks} = \eta_{kks}^0 + \sum_{i<j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jk}\nu_{vs}A_{ijv}$$

$$\xi_{kks} = \xi_{kks}^0 + \sum_{i<j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jk}\nu_{vs}\left(1 - A_{ijv}\right)$$

.

**Proof** In accordance with the principles of variational Bayes, the optimal probability distribution can be formulated as follows:

$$
\begin{aligned}
\log q(\boldsymbol{\alpha}) &\propto \mathrm{E}_{\mathbf{Z},\mathbf{W}}[\log \mathbb{P}(\mathbf{A},\mathbf{Z},\boldsymbol{\alpha},\mathbf{W})]\\
&\propto \mathrm{E}_{\mathbf{Z},\mathbf{W}}[\log p(\mathbf{A}\mid \mathbf{Z},\mathbf{W},\boldsymbol{\alpha})] + \log\mathbb{P}(\boldsymbol{\alpha})\\
&= \sum_{i<j}^{N}\sum_{k,l}^{K}\sum_{v}^{V}\sum_{s}^{Q}\tau_{ik}\tau_{jl}\nu_{vs}\left(A_{ijv}\log(\alpha_{kls}) + (1 - A_{ijv})\log\left(1 - \alpha_{kls}\right)\right)\\
&\quad + \sum_{k\le l}^{K}\sum_{s}^{Q}\left(\left(\eta_{kls}^0 - 1\right)\log(\alpha_{kls}) + \left(\xi_{kls}^0 - 1\right)\log\left(1 - \alpha_{kls}\right)\right)\\
&= \sum_{k<l}^{K}\sum_{i\ne j}^{N}\sum_{v}^{V}\sum_{s}^{Q}\tau_{ik}\tau_{jl}\nu_{vs}\left(A_{ijv}\log(\alpha_{kls}) + (1 - A_{ijv})\log\left(1 - \alpha_{kls}\right)\right)\\
&\quad + \sum_{k=1}^{K}\sum_{i<j}^{N}\sum_{v}^{V}\sum_{s=1}^{Q}\tau_{ik}\tau_{jk}\nu_{vs}\left(A_{ijv}\log(\alpha_{kks}) + (1 - A_{ijv})\log\left(1 - \alpha_{kks}\right)\right)\\
&\quad + \sum_{k\le l}^{K}\sum_{s}^{Q}\left(\left(\eta_{kls}^0 - 1\right)\log(\alpha_{kls}) + \left(\xi_{kls}^0 - 1\right)\log\left(1 - \alpha_{kls}\right)\right)\\
&= \sum_{k<l}^{K}\sum_{s}^{Q}\left(\eta_{kls}^0 - 1 + \sum_{i\ne j}^{N}\sum_{v}^{V}\tau_{ik}\tau_{jl}\nu_{vs}A_{ijv}\right)\log(\alpha_{kls}) +\\
&\qquad \left(\xi_{kls}^0 - 1 + \sum_{i\ne j}^{N}\sum_{v}^{V}\tau_{ik}\tau_{jl}\nu_{vs}\left(1 - A_{ijv}\right)\right)\log\left(1 - \alpha_{kls}\right)\\
&\quad + \sum_{k=1}^{K}\sum_{s}^{Q}\left(\eta_{kks}^0 - 1 + \sum_{i<j}^{N}\sum_{v}^{V}\tau_{ik}\tau_{jk}\nu_{vs}A_{ijv}\right)\log\alpha_{kks} +\\
&\qquad \left(\xi_{kks}^0 - 1 + \sum_{i<j}^{N}\sum_{v}^{V}\tau_{ik}\tau_{jk}\nu_{vs}\left(1 - A_{ijv}\right)\right)\log\left(1 - \alpha_{kks}\right)
\end{aligned}
$$

.

Therefore,

$$q(\alpha_{kls}) = \mathrm{Beta}(\alpha_{kls};\eta_{kls},\xi_{kls}),$$

if $k \neq l$,

$$\eta_{kls} = \eta_{kls}^0 + \sum_{i \neq j}^{N} \sum_{v}^{V} \tau_{ik} \tau_{jl} \nu_{vs} A_{ijv}$$
$$\xi_{kls} = \xi_{kls}^0 + \sum_{i \neq j}^{N} \sum_{v}^{V} \tau_{ik} \tau_{jl} \nu_{vs} \left(1 - A_{ijv}\right)$$

;

otherwise,

$$\eta_{kks} = \eta_{kks}^0 + \sum_{i < j}^{N} \sum_{v}^{V} \tau_{ik} \tau_{jk} \nu_{vs} A_{ijv}$$
$$\xi_{kks} = \xi_{kks}^0 + \sum_{i < j}^{N} \sum_{v}^{V} \tau_{ik} \tau_{jk} \nu_{vs} \left(1 - A_{ijv}\right)$$

.

∎

## Appendix C. Evidence Lower Bound

The lower bound assumes a simplified form after the variational Bayes M-step. It relies solely on the posterior probabilities $\tau_{ik}$ and $\nu_{vs}$ and the normalizing constants of the Dirichlet and Beta distributions.

$$
\begin{aligned}
\mathcal{L}\left(\,q(.)\,\right) = & \log\left\{\frac{\Gamma\left(\sum_{k=1}^{K}\beta_k^0\right)\prod_{k=1}^{K}\Gamma\left(\beta_k\right)}{\Gamma\left(\sum_{k=1}^{K}\beta_k\right)\prod_{k=1}^{K}\Gamma\left(\beta_k^0\right)}\right\} + \log\left\{\frac{\Gamma\left(\sum_{s=1}^{Q}\theta_s^0\right)\prod_{s=1}^{Q}\Gamma\left(\theta_s\right)}{\Gamma\left(\sum_{s=1}^{Q}\theta_s\right)\prod_{s=1}^{Q}\Gamma\left(\theta_s^0\right)}\right\} \\
& + \sum_{k\leq l}^{K}\sum_{s=1}^{Q}\log\left\{\frac{\Gamma\left(\eta_{kls}^0+\xi_{kls}^0\right)\Gamma\left(\eta_{kls}\right)\Gamma\left(\xi_{kls}\right)}{\Gamma\left(\eta_{kls}+\xi_{kls}\right)\Gamma\left(\eta_{kls}^0\right)\Gamma\left(\xi_{kls}^0\right)}\right\} \\
& - \sum_{i}^{N}\sum_{k}^{K}\tau_{ik}\log\tau_{ik} \; - \sum_{v}^{V}\sum_{s}^{Q}\nu_{vs}\log\nu_{vs}
\end{aligned}
$$

**Proof** The lower bound can be expressed as:

$$
\begin{aligned}
\mathcal{L}\left(q(.)\right) = & \sum_{\mathbf{Z}}\sum_{\mathbf{W}}\int\int\int q(\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\pi},\boldsymbol{\rho})\log\frac{\mathbb{P}\left(\mathbf{A},\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\pi},\boldsymbol{\rho}\right)}{q(\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\pi},\boldsymbol{\rho})}\,d\boldsymbol{\alpha}\,d\boldsymbol{\pi}\,d\boldsymbol{\rho} \\
= & \;\mathbb{E}_{\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\rho},\boldsymbol{\pi}}[\log\mathbb{P}(\mathbf{A},\mathbf{Z},\boldsymbol{\alpha},\mathbf{W},\boldsymbol{\rho},\boldsymbol{\pi})] - \mathbb{E}_{\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\rho},\boldsymbol{\pi}}[\log q(\mathbf{Z},\boldsymbol{\alpha},\mathbf{W},\boldsymbol{\rho},\boldsymbol{\pi})]
\end{aligned}
$$

We can decompose the following terms as:

$$
\begin{aligned}
\mathbb{E}_{\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\rho},\boldsymbol{\pi}}[\log p(\mathbf{A},\mathbf{Z},\boldsymbol{\alpha},\mathbf{W},\boldsymbol{\rho},\boldsymbol{\pi})] = & \;\mathbb{E}_{\mathbf{Z},\mathbf{W},\boldsymbol{\alpha}}[\log\mathbb{P}(\mathbf{A}\mid\mathbf{Z},\mathbf{W},\boldsymbol{\alpha})] + \mathbb{E}_{\boldsymbol{\alpha}}[\log p(\boldsymbol{\alpha})] \\
& + \mathbb{E}_{\mathbf{Z},\boldsymbol{\pi}}[\log p(\mathbf{Z}\mid\boldsymbol{\pi})] + \mathbb{E}_{\boldsymbol{\pi}}[\log p(\boldsymbol{\pi})] \\
& + \mathbb{E}_{\mathbf{W},\boldsymbol{\rho}}[\log p(\mathbf{W}\mid\boldsymbol{\rho})] + \mathbb{E}_{\boldsymbol{\rho}}[\log p(\boldsymbol{\rho})]
\end{aligned}
,
$$

and

$$
\begin{aligned}
\mathbb{E}_{\mathbf{Z},\mathbf{W},\boldsymbol{\alpha},\boldsymbol{\rho},\boldsymbol{\pi}}[\log q(\mathbf{Z},\boldsymbol{\alpha},\mathbf{W},\boldsymbol{\rho},\boldsymbol{\pi})] = & \;\mathbb{E}_{\mathbf{Z}}[\log q(\mathbf{Z})] + \mathbb{E}_{\boldsymbol{\pi}}[\log q(\boldsymbol{\pi})] \\
& + \mathbb{E}_{\mathbf{Z}}[\log q(\mathbf{W})] + \mathbb{E}_{\boldsymbol{\rho}}[\log q(\boldsymbol{\rho})] \\
& + \mathbb{E}_{\boldsymbol{\alpha}}[\log q(\boldsymbol{\alpha})]
\end{aligned}
.
$$

Now, the next step involves developing each of these terms and simplifying them as extensively as possible.

$$
\begin{aligned}
\mathbb{E}_{\mathbf{Z},\mathbf{W},\boldsymbol{\alpha}}[\log\mathbb{P}(\mathbf{A}\mid\mathbf{Z},\mathbf{W},\boldsymbol{\alpha})] + \mathbb{E}_{\boldsymbol{\alpha}}[\log\mathbb{P}(\boldsymbol{\alpha})] = & \sum_{i<j}^{N}\sum_{k,l}^{K}\sum_{v}^{V}\sum_{s}^{Q}\tau_{ik}\tau_{jl}\nu_{vs}\Big\{A_{ijv}\Big(\psi(\eta_{kls})-\psi(\xi_{kls})\Big)+\psi(\xi_{kls}) \\
& -\psi(\eta_{kls}+\xi_{kls})\Big\} + \sum_{k\leq l}^{K}\sum_{s}^{Q}\Big\{\log\Gamma(\eta_{kls}^0+\xi_{kls}^0)-\log\Gamma(\eta_{kls}^0) \\
& -\log\Gamma(\xi_{kls}^0)+\left(\eta_{kls}^0-1\right)\left(\psi(\eta_{kls})-\psi(\xi_{kls}+\eta_{kls})\right)+ \\
& \left(\xi_{kls}^0-1\right)\left(\psi(\xi_{kls})-\psi(\eta_{kls}+\xi_{kls})\right)\Big\}
\end{aligned}
$$

$$\mathbb{E}_{\mathbf{Z},\boldsymbol{\pi}}[\log p(\mathbf{Z} \mid \boldsymbol{\pi})] + \mathbb{E}_{\boldsymbol{\pi}}[\log p(\boldsymbol{\pi})] = \sum_i^N \sum_k^K \tau_{ik} \left( \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) \right)$$
$$+ \log \Gamma(\sum_{k'} \beta_{k'}^0) - \log \left( \sum_{k'} \Gamma(\beta_{k'}^0) \right) + \sum_{k=1}^K \left( \beta_k^0 - 1 \right) \left( \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) \right)$$

$$\mathbb{E}_{\mathbf{W},\boldsymbol{\rho}}[\log p(\mathbf{W} \mid \boldsymbol{\rho})] + \mathbb{E}_{\boldsymbol{\rho}}[\log p(\boldsymbol{\rho})] = \sum_v^V \sum_s^Q \nu_{vs} \left( \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) \right)$$
$$+ \log \Gamma(\sum_{s'} \theta_{s'}^0) - \log \left( \sum_{s'} \Gamma(\theta_{s'}^0) \right) + \sum_{s=1}^Q \left( \theta_s^0 - 1 \right) \left( \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) \right)$$

$$\mathbb{E}_{\mathbf{Z}}[\log q(\mathbf{Z})] + \mathbb{E}_{\boldsymbol{\pi}}[\log q(\boldsymbol{\pi})] = \sum_i^N \sum_k^K \tau_{ik} \log \tau_{ik}$$
$$+ \log \Gamma(\sum_{k'} \beta_{k'}) - \log \left( \sum_{k'} \Gamma(\beta_{k'}) \right) + \sum_{k=1}^K \left( \beta_k - 1 \right) \left( \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) \right)$$

$$\mathbb{E}_{\mathbf{Z}}[\log q(\mathbf{W})] + \mathbb{E}_{\boldsymbol{\rho}}[\log q(\boldsymbol{\rho})] = \sum_v^V \sum_s^Q \nu_{vs} \log \nu_{vs}$$
$$+ \log \Gamma(\sum_{s'} \theta_{s'}) - \log \left( \sum_{s'} \Gamma(\theta_{s'}) \right) + \sum_{s=1}^Q \left( \theta_s - 1 \right) \left( \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) \right)$$

$$\mathbb{E}_{\boldsymbol{\alpha}}[\log q(\boldsymbol{\alpha})] = \sum_{k \leq l}^K \sum_s^Q \left\{ \log \Gamma(\eta_{kls} + \xi_{kls}) - \log \Gamma(\eta_{kls}) - \log \Gamma(\xi_{kls}) \right.$$
$$\left. + (\eta_{kls} - 1)(\psi(\eta_{kls}) - \psi(\xi_{kls} + \eta_{kls})) + (\xi_{kls} - 1)(\psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls})) \right\}$$

Now that all the terms have been developed, it's just a matter of grouping them together, to obtain the ELBO below.

$$
\begin{aligned}
\mathcal{L}\left(q(.)\right) = & \sum_{k<l}^{K} \sum_{s}^{Q} \left( \eta_{kls}^0 + \left( \sum_{i\neq j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jl}\nu_{vs}A_{ijv} \right) - \eta_{kls} \right) \left( \psi(\eta_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right) \\
& + \sum_{k=1}^{K} \sum_{s}^{Q} \left( \eta_{kks}^0 + \left( \sum_{i<j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jk}\nu_{vs}A_{ijv} \right) - \eta_{kks} \right) \left( \psi(\eta_{kks}) - \psi(\eta_{kks} + \xi_{kks}) \right) \\
& + \sum_{k<l}^{K} \sum_{s}^{Q} \left( \xi_{kls}^0 + \left( \sum_{i\neq j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jl}\nu_{vs}(1 - A_{ijv}) \right) - \eta_{kls} \right) \left( \psi(\xi_{kls}) - \psi(\eta_{kls} + \xi_{kls}) \right) \\
& + \sum_{k=1}^{K} \sum_{s}^{Q} \left( \xi_{kks}^0 + \left( \sum_{i<j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jk}\nu_{vs}(1 - A_{ijv}) \right) - \xi_{kks} \right) \left( \psi(\xi_{kks}) - \psi(\eta_{kks} + \xi_{kks}) \right) \\
& + \sum_{k=1}^{K} \left( \beta_k^0 + \sum_{i=1}^{N} \tau_{ik} - \beta_k \right) \left( \psi(\beta_k) - \psi(\sum_{k'} \beta_{k'}) \right) \\
& + \sum_{q=1}^{Q} \left( \theta_s^0 + \sum_{v=1}^{V} \nu_{vs} - \theta_s \right) \left( \psi(\theta_s) - \psi(\sum_{s'} \theta_{s'}) \right) \\
& - \sum_{i}^{N} \sum_{k}^{K} \tau_{ik} \log \tau_{ik} - \sum_{v}^{V} \sum_{s}^{Q} \nu_{vs} \log \nu_{vs} \\
& + \log \left\{ \frac{\Gamma\left(\sum_{k=1}^{K} \beta_k^0\right) \prod_{k=1}^{K} \Gamma(\beta_k)}{\Gamma\left(\sum_{k=1}^{K} \beta_k\right) \prod_{k=1}^{K} \Gamma(\beta_k^0)} \right\} + \log \left\{ \frac{\Gamma\left(\sum_{s=1}^{Q} \theta_s^0\right) \prod_{s=1}^{Q} \Gamma(\theta_s)}{\Gamma\left(\sum_{s=1}^{Q} \theta_s\right) \prod_{s=1}^{Q} \Gamma(\theta_s^0)} \right\} \\
& + \sum_{k\leq l}^{K} \sum_{s=1}^{Q} \log \left\{ \frac{\Gamma\left(\eta_{kls}^0 + \xi_{klq}^0\right) \Gamma(\eta_{kls}) \Gamma(\xi_{kls})}{\Gamma(\eta_{kls} + \xi_{kls}) \Gamma(\eta_{kls}^0) \Gamma(\xi_{kls}^0)} \right\}
\end{aligned}
$$

However, by definition of the parameters, we have many terms that cancel each other out:

- $\eta_{kls} = \eta_{kls}^0 + \left( \sum_{i\neq j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jl}\nu_{vs}A_{ijv} \right)$

- $\eta_{kks} = \eta_{kks}^0 + \left( \sum_{i<j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jk}\nu_{vs}A_{ijv} \right)$

- $\eta_{kls} = \xi_{kls}^0 + \left( \sum_{i\neq j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jl}\nu_{vs}(1 - A_{ijv}) \right)$

- $\xi_{kks} = \xi_{kks}^0 + \left( \sum_{i<j}^{N} \sum_{v}^{V} \tau_{ik}\tau_{jk}\nu_{vs}(1 - A_{ijv}) \right)$

- $\beta_k = \beta_k^0 + \sum_{i=1}^{N} \tau_{ik}$

- $\theta_s = \theta_s^0 + \sum_{v=1}^{V} \nu_{vs}$

Hence:

$$\mathcal{L}\left(q(.)\right) = \log\left\{\frac{\Gamma\left(\sum_{k=1}^{K}\beta_k^0\right)\prod_{k=1}^{K}\Gamma\left(\beta_k\right)}{\Gamma\left(\sum_{k=1}^{K}\beta_k\right)\prod_{k=1}^{K}\Gamma\left(\beta_k^0\right)}\right\} + \log\left\{\frac{\Gamma\left(\sum_{s=1}^{Q}\theta_s^0\right)\prod_{s=1}^{Q}\Gamma\left(\theta_s\right)}{\Gamma\left(\sum_{s=1}^{Q}\theta_s\right)\prod_{s=1}^{Q}\Gamma\left(\theta_s^0\right)}\right\}$$

$$+ \sum_{k\leq l}^{K}\sum_{s=1}^{Q}\log\left\{\frac{\Gamma\left(\eta_{kls}^0 + \xi_{klq}^0\right)\Gamma\left(\eta_{kls}\right)\Gamma\left(\xi_{kls}\right)}{\Gamma\left(\eta_{kls} + \xi_{kls}\right)\Gamma\left(\eta_{kls}^0\right)\Gamma\left(\xi_{kls}^0\right)}\right\}$$

$$- \sum_{i}^{N}\sum_{k}^{K}\tau_{ik}\log\tau_{ik} \ - \sum_{v}^{V}\sum_{s}^{Q}\nu_{vs}\log\nu_{vs}$$

∎