

1. Introduction / Abstract:

This paper is motivated ^{by} a problem of network clustering where each node in a network has a covariate or covariates. That is, given some adjacency matrix x ($n \times n$ $x_{ij} \in \{1, 2, \dots, n\}$) where $\text{diag}(x) = 0$ and x is symmetric and some node specific features as a set $y = \{y_1, y_2, \dots, y_n\}$ where each y_i is p -dimensional vectors of features, How can we use this information to find clusters?

They propose to use classical Stochastic Blockmodel (SBM) and extend it for this case. They call this new extension the Mix Net model.

2. A mixture of Network with covariates

The Mix Net assumes that the covariates y and edges x are independent given the node clusters. Or in other words that the connectivity profile of a given node and its features can be explained by its cluster.

To comment on this further, we introduce some notation.

- x is $n \times n$ matrix $i, j \in \{1, 2, \dots, n\}$ nodes thus $x_{ij} = 1$ or $x_{ij} = 0$.
- X is $n \times n$ random matrix
- Z is $n \times Q$ matrix of $i \in \{1, 2, \dots, n\}$ $q \in \{1, 2, \dots, Q\}$ clusters. Thus $Z_1 = (0, 1, 0, 0)$ means there was $Q=4$ and node 1 is located in the 2nd class.
- Z is random variable.
- y is $n \times p$ matrix x of features $i \in \{1, 2, \dots, n\}$ $p \in \{1, \dots, p\}$.
- y is random variable.

Distributions:

- $x_{ij} | Z_{iq}=1, Z_{jq}=1 \sim \text{Bernoulli}(\pi_{qk})$ small letters.
- $f(x|Z; \pi) = \prod_{i=1}^n \prod_{j \neq i}^n \prod_{q=1}^Q \prod_{k=1}^Q f(x_{ij} | Z_{iq}=1, Z_{jq}=1)^{\frac{Z_{iq} Z_{jq}}{2}}$ (EQ. 1)
- $Z_i \sim \text{Categorical}(d)$, where d is a vector of length Q .
- $(y_i | Z_{iq}=1) \sim \text{MVN}(\mu_q, \Sigma_q)$ $f(y_i | Z_{iq}) = \prod_{l=1}^p \prod_{g=1}^a f(y_{il} | Z_{iq})$ (EQ. 3)
for each group variance component
- $\mu_q = (\mu_q^{(1)}, \dots, \mu_q^{(p)})$ and $\Sigma_q = \sigma_q^2 I$ $p \times p$ matrix.

The latter assumption $\Sigma_a = b^2 I$ means that all features within a cluster z have a common variance and moreover that this variance is constant ^{within a} clusters. Thus, it is important to standardise node features before the analysis. (3)

Model's likelihood.

$$\begin{aligned} \log f(x, y, z; \pi, \mu, b^2, d) &= \log f(x, y | z; \pi, \mu, b^2) f(z; d) \\ &= \log f(x | z; \pi) + \log f(y | z; \mu, b^2) \\ &\quad + \log f(z; d) \quad \mathbb{E}(\mathcal{Y}) \end{aligned}$$

In particular, they are interested in optimising marginal log-likelihood:

$$\log f(x, y; \pi, \mu, b^2, d) = \log \sum_{z \in S} f(x, y, z; \pi, \mu, b^2, d)$$

which is not tractable. So they propose to use variational approximation. In general, variational approximation is understood in the following way:

$$\log f(x, y; \pi, \mu, b^2, d) = \log \sum_{z \in S} f(x, y, z; \pi, \mu, b^2, d) \frac{f^*(z; \pi)}{f^*(z; \pi)}$$

$$= \log \mathbb{E} \left[\frac{f(x, y, z; \pi, \mu, b^2, d)}{f^*(z; \pi)} \right]$$

$$\geq \mathbb{E} \left[\log \frac{f(x, y, z; \pi, \mu, b^2, d)}{f^*(z; \pi)} \right] \quad \text{by Jensen's inequality.}$$

$$= \mathbb{E}_{f^*} \left[\log f(x, y, z; \pi, \mu, b^2, d) \right] - \mathbb{E}_{f^*} \left[\log f^*(z; \pi) \right]$$

In particular, the lower bound (RHS) is closely related to the Kullback Leibler divergence $KL(\cdot)$ of $f^*(z; \pi)$ to $f(z | x, y; \pi, \mu, b^2, d)$.

$$KL[f^*(z; \pi) | f(z | x, y; \pi, \mu, b^2, d)] = \log f(x, y; \pi, \mu, b^2, d) - \left\{ \mathbb{E} \left[\log f(x, y, z; \pi, \mu, b^2, d) \right] - \mathbb{E} \left[\log f^*(z; \pi) \right] \right\}$$

Rearranging

$$\log f(x, y; \pi, \mu, b^2, d) = KL[\cdot] = \mathbb{E} \left[\log f(x, y, z; \pi, \mu, b^2, d) \right] - \mathbb{E} \left[\log f^*(z; \pi) \right]$$

Thus the lower bound is attained when $KL[\cdot] = 0$ that is when $f^*(z; \pi)$ coincides with $f(z | x, y; \pi, \mu, b^2, d)$. Thus, the goal is for a family of densities described by variational parameter π , $[n \times Q$ matrix of probabilities] we can approximate marginal likelihood by its lower bound.

Noting the lower bound as:

$$J(f^*(z; \tau); \pi, d, \mu, b^2) = E[\log f(x|y, z; d, \pi, \mu, b^2)] - E[\log f^*(z; \tau)] \quad (5)$$

Now the point estimating equations for d, π, b^2, μ, τ are found by maximizing $J(\cdot)$ w.r.t to the model's parameters.

$$J(\cdot) = E[\log f(x|y|z; \pi, \mu, b^2) f(z; d)] - E[\log f^*(z; \tau)]$$

$$= E[\underbrace{\log f(x|z; \pi)}_{(EQ. 1)} \underbrace{f(y|z; \mu, b^2)}_{(EQ. 3)} \underbrace{f(z; d)}_{(EQ. 2)}] - E[\underbrace{\log f^*(z; \tau)}_{(EQ. 5)}]$$

$$= E[\log \prod_{i=1}^n \prod_{j=1}^n \prod_{q=1}^a f(x_{ij} | z_{iq}, z_{iq}, \pi_{jq})] + E[\log \prod_{i=1}^n \prod_{q=1}^a f(y_i | z_{iq})] \\ + E[\log \prod_{i=1}^n \prod_{q=1}^a d_{iq}^{z_{iq}}] - E[\log \prod_{i=1}^n \prod_{q=1}^a \tau_{iq}^{z_{iq}}]$$

After taking expt. $E(z_{iq} z_{il}) = \tau_{iq} \tau_{il}$ and $E(z_{iq}) = \tau_{iq}$

$$J(\cdot) = \sum_{i=1}^n \sum_{j=1}^n \sum_{q=1}^a \tau_{iq} \tau_{jl} \log f(x_{ij} | z_{iq}, z_{iq}; \pi_{jq}) + \sum_{i=1}^n \sum_{q=1}^a \tau_{iq} \log f(y_i | z_{iq}) \\ + \sum_{i=1}^n \sum_{q=1}^a \tau_{iq} \log d_{iq} - \sum_{i=1}^n \sum_{q=1}^a \tau_{iq} \log \tau_{iq}$$

For example:

$$\frac{\partial}{\partial d_{iq}} J(\cdot) = \frac{\partial}{\partial d_{iq}} \left[\sum_{i=1}^n \sum_{q=1}^a \tau_{iq} \log d_{iq} \right] \Rightarrow \\ \Rightarrow \sum_{i=1}^n \tau_{iq} \frac{1}{d_{iq}} + \frac{\partial}{\partial d_{iq}} \left[\lambda \sum_{i=1}^n \sum_{q=1}^a d_{iq} - 1 \right] = 0$$

$$\frac{1}{d_{iq}} \sum_{i=1}^n \tau_{iq} + \lambda = 0$$

$$\frac{1}{d_{iq}} \sum_{i=1}^n \tau_{iq} = \lambda$$

$$d_{iq} = \frac{\sum_{i=1}^n \tau_{iq}}{\lambda}$$

But we know that

$$\sum_{q=1}^a d_{iq} = \sum_{q=1}^a \frac{1}{\lambda} \sum_{i=1}^n \tau_{iq}$$

$$1 = \frac{1}{\lambda} \sum_{i=1}^n \sum_{q=1}^a \tau_{iq}$$

$$1 = \frac{1}{\lambda} \sum_{i=1}^n \sum_{q=1}^a \tau_{iq}$$

$$1 = \frac{1}{\lambda} \sum_{i=1}^n 1$$

$$1 = \frac{1}{\lambda} \cdot n$$

$$\lambda = \frac{n}{1} = n$$

$$d_{iq} = \frac{1}{n} \sum_{i=1}^n \tau_{iq}$$

Similarly, the updating equations are:

$$\hat{\pi}_{ij} = \frac{\sum_{i \neq j} t_{ij} t_{je} x_{ij}}{\sum_{i \neq j} t_{ij} t_{je}}$$

$$\hat{\mu}_a = \frac{\sum_{i=1}^n \hat{\pi}_{ij} y_i}{\sum_{i=1}^n \hat{\pi}_{ij}}$$

$$\hat{\sigma}_a^2 = \frac{\sum_{i=1}^n \sum_{j=1}^q t_{ij} (y_i - \hat{\mu}_a)^T (y_i - \hat{\mu}_a)}{P \sum_{i=1}^n \sum_{j=1}^q \hat{\pi}_{ij}} = n \quad \left(\text{Note, eq (a) contains 5 hyperparameters } \sigma_a^2 = ? \text{ should be } \mu_a \right)$$

$$\hat{\pi}_{ij} \propto \frac{\hat{\sigma}_a^2}{\hat{\sigma}_a^2} \prod_{i \neq j} \prod_{\ell=1}^q \left[\frac{x_{ij}^{\pi_{ij}} (1 - \pi_{ij})^{1 - x_{ij}}}{\pi_{ij}^{\pi_{ij}} (1 - \pi_{ij})^{1 - \pi_{ij}}} \right] \prod_{p=1}^P \exp \left\{ - \frac{(y_i^{(p)} - \hat{\mu}_a^{(p)})^2}{2 \hat{\sigma}_a^2} \right\}$$

Different to Eq (8), as the parameters we try to estimate are never dropped. b b 1

Model Selection: ICL

So far we assumed that Q was fixed. Here, we will discuss criterion for its estimation. Here, we assume that a model with Q blocks (clusters) is random variable and that all m_a 's are equally likely. We base this criterion on the integrated classification likelihood, that is:

$$f(x, y, z | m_a) = f(x | z, m_a) f(y | z, m_a) f(z | m_a)$$

Also, we will assume that all priors are separable.

$$f(x | z, m_a) = \int_{\Pi_{[Q \times Q]}} f(x | z, \Pi, m_a) f(\Pi | m_a) d\Pi$$

It is assumed that the conditions of Laplace integral approximation are satisfied, that is that the contribution to the integral is the greatest in the neighborhood of MLE estimate $\hat{\Pi}$ and that the influence of prior $f(\Pi | m_a)$ is negligible.

$$\log f(x | z, m_a) \approx \log f(x | z, \hat{\Pi}, m_a) - \frac{1}{2} \frac{Q(Q-1)}{2} \log \binom{n}{2}$$

$$\log f(y | z, m_a) \approx \log f(y | z, \hat{\mu}, \hat{\sigma}^2, m_a) - \frac{Q}{2} (P+1) \log(nP)$$

$$\log f(z | m_a) \approx \log f(z | m_a, \hat{\pi}) - \frac{Q(Q-1)}{2} \log(n)$$

Combining these we obtain the IC (ma).

$$IC(m_a) = \log f(x, y, z | m_a, \hat{\lambda}, \hat{\pi}, \hat{M}, \hat{\sigma}^2) - \frac{Q}{2} (P+1) \log(nP) - \frac{1}{2} \frac{Q(Q-1)}{2} \log\left(\frac{n(n-1)}{2}\right) - \frac{(Q-1)}{2} \log(n).$$

9

This construction is \neq to the one given in the paper, as the criterion they propose rewards for complexity.

Experiments. (This is summarized in Table 1)

They consider PIs of the form $\Sigma = \begin{bmatrix} \lambda & \epsilon & \dots & \epsilon \\ & \ddots & & \\ & & \lambda & \epsilon \\ & & & \ddots \end{bmatrix}$ $d_i = \frac{1}{a}$ and $n=180$

Here 2 parameters model the distance between $Q \times Q$ within-cluster and between-cluster Bernoulli rates. Similarly, they consider $(M_2$ and $M_4)$ and this "distance" can also capture some difficult-easy cases. They consider 43 different scenarios for these models and the simulations per ~~model~~ case is 20. To evaluate if the clusters assignments make sense they use ARI (Adjusted Rand Index) which evaluates similarity on $[0,1]$ scale +1 being identical and 0 being totally different. For each, simulation alg. is restarted 10 times. This

Alternative algorithms.

- ① Spectral Multiple View Learning (SMVL)
- ② Hidden Markov Random Fields. (Ambrus et al. 1997), (not sure if these 2 can model for vertex features?)

Results.

In the Figure 2 (a) (b) (c) when there is strong evidence of modular structure in the network (i.e. $d(\lambda, \epsilon) > 0$, $d(M_2, M_4) > 0$) then the CohMix algorithm performs better than SMVL and HMRF algorithms. This is of course measured in terms of ARI scores. What are called Rand Errors?