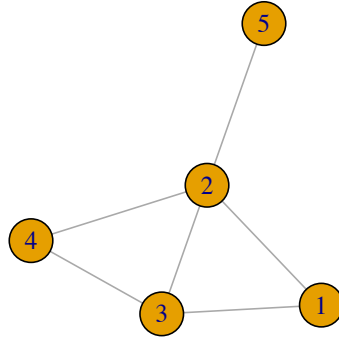


Mock exam for Statistical analysis of network data - solutions

1. (a)



(b)

$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix}$$

(c)

$$I = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{pmatrix}$$

(d) By definition, $d_i = \sum_j A_{ij}$. Substituting, $d_1 = 2$, $d_2 = 4$, $d_3 = 3$, $d_4 = 2$, $d_5 = 2$.

(e) $N_{C_3} = \frac{1}{6} \sum_{i,j,k} A_{ij} A_{jk} A_{ki} = 2$.

2. We define

$$d_i = \sum_{j \neq i} A_{ij}. \tag{1}$$

(a) To calculate the moments of this random variable we use the law of iterated expectation

$$\begin{aligned} E\{d_i\} &= E_{\xi} E_{A|\xi} \{d_i\} = \sum_{j \neq i} E_{\xi} g(\xi_i) g(\xi_j) \\ &= (n-1) E^2\{g(\xi_i)\} = (n-1) \left(\int_0^1 g(x) \cdot 1 dx \right)^2 = (n-1) \|g\|_1^2. \end{aligned}$$

We define the function p -norm to be

$$\|g\|_p^p = \int_0^1 g^p(x) dx.$$

(b) To compute the variance we use the law of total variance as seen in class:

$$\begin{aligned}
\text{Var}\{d_i\} &= \mathbb{E}_\xi \text{Var}_{A|\xi}\{d_i\} + \text{Var}_\xi \mathbb{E}_{A|\xi}\{d_i\} \\
&= \mathbb{E}_\xi \sum_{j \neq i} g(\xi_i)g(\xi_j)(1 - g(\xi_i)g(\xi_j)) + \text{Var}_\xi \sum_{j \neq i} g(\xi_i)g(\xi_j) \\
&= (n-1)\{\|g\|_1^2 - \|g\|_2^4\} + \mathbb{E} \sum_{j \neq i} \sum_{k \neq i} g(\xi_i)g(\xi_j)g(\xi_k)g(\xi_i) - (n-1)^2 \|g\|_1^4 \\
&= (n-1)\{\|g\|_1^2 - \|g\|_2^4\} + \mathbb{E} \sum_{j \neq i} \sum_{k \neq i} g(\xi_i)g(\xi_j)g(\xi_k)g(\xi_i) - (n-1)^2 \|g\|_1^4 \\
&= (n-1)\{\|g\|_1^2 - \|g\|_2^4\} + \mathbb{E} \sum_{j \neq i} \sum_{k \neq i, j} g^2(\xi_i)g(\xi_j)g(\xi_k) + \mathbb{E} \sum_{j \neq i} g^2(\xi_i)g^2(\xi_j) \\
&\quad - (n-1)^2 \|g\|_1^4 \\
&= (n-1)\{\|g\|_1^2 - \|g\|_2^4\} + (n-1)(n-2)\|g\|_2^2 \|g\|_1^2 + (n-1)\|g\|_2^4 - (n-1)^2 \|g\|_1^4 \\
&= (n-1)\|g\|_1^2 + (n-1)(n-2)\|g\|_2^2 \|g\|_1^2 - (n-1)^2 \|g\|_1^4.
\end{aligned}$$

For the special case of an Erdos Renyi network this becomes

$$\begin{aligned}
\text{Var}\{d_i\} &= (n-1)\rho_n + (n-1)(n-2)\rho_n^2 - (n-1)^2 \rho_n^2 \\
&= (n-1)\rho_n + (n-1)\{n-2 - (n-1)\}\rho_n^2 = (n-1)\rho_n + (n-1)(-1)\rho_n^2.
\end{aligned}$$

We can also note that if A is an Erdos–Renyi graph with edge probability ρ_n then $\|g\|_1 = \sqrt{\rho_n}$, then $d_i = \sum_{j \neq i} A_{ij} \sim \text{Bin}\{(n-1), \rho_n\}$ which has mean $(n-1)\rho_n$ and variance $(n-1)\rho_n(1 - \rho_n)$, which matches.

We can also compute this directly. We already have $\mathbb{E}\{d_i\}$ above. Then it follows

$$\begin{aligned}
\mathbb{E}\{d_i^2\} &= \mathbb{E} \sum_{j \neq i} \sum_{k \neq i} A_{ij} A_{ik} \\
&= \mathbb{E} \sum_{j \neq i} \sum_{k \neq i, j} A_{ij} A_{ik} + \mathbb{E} \sum_{j \neq i} A_{ij} \\
&= \|g\|_2^2 \|g\|_1^2 (n-1)(n-2) + \|g\|_1^2 (n-1).
\end{aligned} \tag{2}$$

We now note that:

$$\begin{aligned}
\text{Var}\{d_i\} &= \mathbb{E}\{d_i^2\} - \mathbb{E}^2\{d_i\} \\
&= \|g\|_2^2 \|g\|_1^2 (n-1)(n-2) + \|g\|_1^2 (n-1) - (n-1)^2 \|g\|_1^4.
\end{aligned} \tag{3}$$

This gives the same result as the law of total variance.

(c) We can calculate this by brute force for $i \neq j$:

$$\begin{aligned}
\text{Cov}\{d_i, d_j\} &= \mathbb{E}\{d_i d_j\} - \mathbb{E}\{d_i\}\mathbb{E}\{d_j\} \\
&= \mathbb{E}\left\{\sum_{k \neq i} A_{ik} \sum_{l \neq j} A_{jl}\right\} - \mathbb{E}\{d_i\}\mathbb{E}\{d_j\} \\
&= \mathbb{E}\left\{\sum_{k \neq i, j, l} A_{ik} \sum_{l \neq j, i} A_{jl}\right\} + \mathbb{E}\left\{\sum_{k=j} \sum_{l=i} A_{ik} A_{jl}\right\} + \mathbb{E}\left\{\sum_{l \neq j} A_{ij} A_{jl}\right\} + \mathbb{E}\left\{\sum_{l \neq j} A_{il} A_{jl}\right\} \\
&\quad + \mathbb{E}\left\{\sum_{l \neq j} A_{il} A_{ji}\right\} - (n-1)^2 \|g\|_1^4 \\
&= \sum_{k \neq i, j, l} \sum_{l \neq j, i} \|g\|_1^4 + \|g\|_2^2 + 3(n-2)\|g\|_2^2 \|g\|_1^2 - (n-1)^2 \|g\|_1^4 \\
&= (n-3)(n-2)\|g\|_1^4 - (n-1)^2 \|g\|_1^4 + \|g\|_2^2 + 3(n-2)\|g\|_2^2 \|g\|_1^2 \\
&= (n^2 - 5n + 6 - n^2 + 2n - 1)\|g\|_1^4 + \|g\|_2^2 + 3(n-2)\|g\|_2^2 \|g\|_1^2 \\
&= (-3n + 5)\|g\|_1^4 + 3(n-2)\|g\|_2^2 \|g\|_1^2 + \|g\|_2^2.
\end{aligned} \tag{4}$$

Therefore in the special case of an Erdos–Renyi graph we get as $\|g\|_p = \sqrt{\rho_n}$ for $p = 1, 2, 3, \dots$ the simple form of

$$\text{Cov}\{d_i, d_j\} = \rho_n(1 - \rho_n),$$

and this concurs with our expectation, as one edge will be in common between the two sums.

Part (ii):

- (a) A graph model is finitely exchangeable if for all permutations $\pi : [n] \rightarrow [n]$, $\mathbf{A} \stackrel{d}{=} \mathbf{A}^{(\pi)}$, where $\mathbf{A}^{(\pi)}$ is the permuted adjacency matrix with elements $A_{ij}^{(\pi)} = A_{\pi_i, \pi_j}$.
- (b) The joint probability is

$$\Pr\{\mathbf{A} = \mathbf{a}\} = \Pr\{\mathbf{A} = \mathbf{a} \mid \boldsymbol{\xi}\} \Pr\{\boldsymbol{\xi}\} = \prod_{i < j} g(\xi_i) g(\xi_j) \Pr\{\xi_i\} \Pr\{\xi_j\},$$

which is obviously invariant to permutations of the indices.

- 3. (i)

$$D = \begin{pmatrix} 0 & 2 & 1 & 2 & 2 \\ 2 & 0 & 1 & 2 & 2 \\ 1 & 1 & 0 & 1 & 1 \\ 2 & 2 & 1 & 0 & 1 \\ 2 & 2 & 1 & 1 & 0 \end{pmatrix}$$

- (ii) Closeness centrality:

$$C_i = \frac{n}{\sum_{j \neq i} D_{ij}} = \frac{5}{\sum_{j \neq i} D_{ij}},$$

from which $C_1 = 5/7$, $C_2 = 5/7$, $C_3 = 5/4$, $C_4 = 5/6$, $C_5 = 5/6$.

- (iii) Harmonic centrality:

$$C_i^{(H)} = \sum_{j \neq i} \frac{1}{D_{ij}},$$

from which $C_1^{(H)} = 2.5$, $C_2^{(H)} = 2.5$, $C_3^{(H)} = 4$, $C_4^{(H)} = 3$, $C_5^{(H)} = 3$.

- 4. (i)

$$\Pr\{\mathbf{A} = \mathbf{a} \mid \mathbf{z}\} = \prod_{i < j} \theta_{z_i z_j}^{a_{ij}} (1 - \theta_{z_i z_j})^{1 - a_{ij}},$$

from which the log-likelihood is

$$\ell(\boldsymbol{\theta}; \mathbf{a}, \mathbf{z}) = \sum_{i < j} \{a_{ij} \log(\theta_{z_i z_j}) + (1 - a_{ij}) \log(1 - \theta_{z_i z_j})\}.$$

- (ii) Described in slides 14-17 of Lecture 8 (Latent space models, April 7).

- (iii) Described in slides 4-5 of Week 6 (Network clustering, March 25).

- 5. (i) The degree-corrected stochastic blockmodel is written as

$$A_{ij} \mid z_i, z_j, \xi_i, \xi_j \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta_{z_i z_j} + \pi(\xi_i)\pi(\xi_j)).$$

From this, using the shorthand $P_{ij} = \pi(\xi_i)\pi(\xi_j)$, the log-likelihood is

$$\ell(\boldsymbol{\theta}; \mathbf{a}, \mathbf{z}, \boldsymbol{\xi}) = \sum_{i < j} \{a_{ij} \log(\theta_{z_i z_j} + P_{ij}) + (1 - a_{ij}) \log(1 - \theta_{z_i z_j} - P_{ij})\}.$$

Assuming all latent variables and the degree structure π known, we can separate out the likelihood component pertaining to each of the blocks:

$$\ell(\theta_{ab}; \mathbf{a}, \mathbf{z}, \boldsymbol{\xi}) = \sum_{i < j} \{a_{ij} \log(\theta_{ab} + P_{ij}(\boldsymbol{\xi})) + (1 - a_{ij}) \log(1 - \theta_{ab} - P_{ij}(\boldsymbol{\xi}))\} I(z_i = a, z_j = b).$$

We can maximize this likelihood with respect to θ_{ab} . The solution will depend on both \mathbf{z} and $\boldsymbol{\xi}$. In principle, we can obtain a global solution by considering the estimator of θ_{ab} for all possible configurations of \mathbf{z} and $\boldsymbol{\xi}$, and find the optimum of the likelihood, but this is very impractical because of the high number of possible cluster membership configurations. Moreover, the presence of the uniform $\boldsymbol{\xi}$ requires using Bayesian methods, marginalising over $\boldsymbol{\xi}$. The EM algorithm may also be used.

- (ii) Introducing the notations $\bar{\theta} = \sum_{a,b} \frac{h_a h_b}{n^2} \theta_{ab}$ and $P = \int_{[0,1]} \pi(x) dx$,

$$\begin{aligned} \mathbb{E}(d_i) &= \sum_j \mathbb{E}_{\boldsymbol{\xi}, \mathbf{z}}[\mathbb{E}(A_{ij} | \boldsymbol{\xi}, \mathbf{z})] = \sum_j \mathbb{E}_{\boldsymbol{\xi}, \mathbf{z}}[\theta_{z_i z_j} + \pi(\xi_i) \pi(\xi_j)] \\ &= \sum_j \mathbb{E}_{\mathbf{z}}(\theta_{z_i z_j}) + \mathbb{E}_{\xi_i}[\pi(\xi_i)] \sum_j \mathbb{E}_{\xi_j}[\pi(\xi_j)] \\ &= \sum_j \bar{\theta} + \sum_j P^2 = (n-1)(\bar{\theta} + P^2). \end{aligned}$$

- (iii) A_{ij} and A_{ji} represent now edges between the same nodes but with opposite direction. Their probability may differ, so the generating model is not necessarily symmetric, neither the interaction matrix θ_{ab} , nor the degree correction. Moreover, even one single node can belong to different clusters when being a donor and when being a receptor, and we might also specify different uniform latent variables for these roles. The model becomes

$$A_{ij} | \boldsymbol{\xi}^{(\text{in})}, \mathbf{z}^{(\text{in})}, \boldsymbol{\xi}^{(\text{out})}, \mathbf{z}^{(\text{out})} \stackrel{\text{iid}}{\sim} \text{Bernoulli}\{\theta_{z_i^{(\text{in})} z_j^{(\text{out})}} + \pi^{(\text{in})}(\xi_i^{(\text{in})}) \pi^{(\text{out})}(\xi_j^{(\text{out})})\}.$$

The corresponding log-likelihood is

$$\begin{aligned} \ell(\theta; \boldsymbol{\xi}^{(\text{in})}, \mathbf{z}^{(\text{in})}, \boldsymbol{\xi}^{(\text{out})}, \mathbf{z}^{(\text{out})}) &= \sum_{i,j} \{a_{ij} \log(\theta_{z_i^{(\text{in})} z_j^{(\text{out})}} + \pi^{(\text{in})}(\xi_i^{(\text{in})}) \pi^{(\text{out})}(\xi_j^{(\text{out})})) \\ &\quad + (1 - a_{ij}) \log(1 - \theta_{z_i^{(\text{in})} z_j^{(\text{out})}} - \pi^{(\text{in})}(\xi_i^{(\text{in})}) \pi^{(\text{out})}(\xi_j^{(\text{out})}))\}. \end{aligned}$$