



# Community Detection in Networks using Probabilistic Graphical Models

## DE424 Mini Project

Jake Love  
25892339

May 12, 2025

### Declaration

I declare that this report and the work presented in it are my own work. I confirm that: this work was done wholly while in candidature for the DE424 module at Stellenbosch University; where I have consulted the published work of others, this is always clearly attributed; where I have quoted from the work of others, the source is always given; and I have acknowledged all main sources of help.

Signed: JGL

Date: 15/05/25

### ECSA Graduate Attributes (GAs)

This mini-project demonstrates the following ECSA Graduate Attributes:

#### GA 1 – Problem Solving

This report demonstrates my ability to identify, formulate, analyze, and solve complex engineering problems creatively and innovatively through: formulating an appropriate PGM-based approach to community detection (Sec. 3); applying mathematical reasoning (Sec. 3, 4); evaluating design options (Sec. 3, 5); and implementing and testing the solution (Sec. 6).

#### GA 2 – Application of Scientific and Engineering Knowledge

This report demonstrates my ability to apply knowledge of mathematics, natural sciences, engineering fundamentals, and engineering specialties through: applying probability theory and inference (Sec. 3, 4); implementing graphical models programmatically (Sec. 4); employing optimization for parameter estimation (Sec. 5); and performing quantitative analysis (Sec. 6).

## 1 Introduction

Community detection, identifying densely connected node groups in networks, is vital for understanding complex systems. This project applies Probabilistic Graphical Models (PGMs), specifically Stochastic Block Models (SBMs), to this task. The objectives are: (1) Design and

implement an SBM-based PGM; (2) Select inference methods; (3) Learn model parameters; (4) Evaluate on synthetic and real-world data. This report details my approach and results.

## 2 Background

A network graph represents entities (nodes) and their interactions (edges). Communities are groups of nodes more densely connected internally than externally. Community detection aims to find these groups from the network structure, often given as an adjacency matrix  $\mathbf{A}$  [2].

Stochastic Block Models (SBMs) are generative models for graphs with communities [1]. My simple SBM assumes  $N$  nodes assigned to one of  $R$  communities. Edges are formed with probability  $w_s$  between nodes in the same community and  $w_d$  between nodes in different communities, usually  $w_s > w_d$ . While SBMs help algorithm development, they may simplify real networks (e.g., non-overlapping communities). PGMs allow model-based inference for community detection using SBMs.

## 3 Objective 1: Model Design and Implementation

### 3.1 Model Selection and Rationale

I implemented a PGM directly modeling the SBM generative process for several reasons:

- It naturally captures the causal relationship between community assignments and edge formation
- This approach allows for efficient inference and parameter estimation with emdw
- The model scales reasonably to larger networks with more communities
- It provides a principled way to determine the optimal number of communities

I explicitly chose to represent communities as discrete random variables rather than binary RVs between node pairs because this:

- Simplifies the model structure and reduces the number of variables
- Makes parameter estimation more straightforward
- Avoids the need to model complex transitivity relationships
- Aligns directly with the SBM generative process

### 3.2 Model Structure

The model includes:

- **Community Variables** ( $C_i$ ): For each node  $i$ , a discrete RV with domain  $\{0, \dots, R-1\}$ , indicating community assignment.
- **Edge Variables** ( $E_{ij}$ ): Observed binary RV for each potential edge  $(i, j)$ , indicating presence ( $E_{ij} = 1$ ) or absence ( $E_{ij} = 0$ ).

The joint distribution is  $P(\mathbf{C}, \mathbf{E}) = P(\mathbf{C})P(\mathbf{E}|\mathbf{C})$ . The overall generative structure involves factors linking community variables ( $C_i, C_j$ ) to the corresponding edge variable ( $E_{ij}$ ).

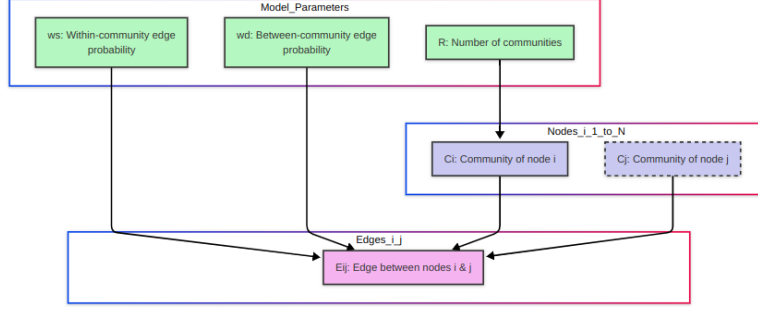


Figure 1: Plate notation representing the generative process of the Stochastic Block Model. This diagram accurately illustrates the hierarchical structure of the model, with community assignments influencing edge probabilities.

### 3.3 Probability Distributions

- **Prior**  $P(C_i)$ : Uniform with small noise to break symmetry,  $P(C_i = r) = 1/R + \varepsilon$
- **Conditional Edge Probability**  $P(E_{ij}|C_i, C_j)$ :

$$P(E_{ij} = 1|C_i = r, C_j = s) = \begin{cases} w_s & \text{if } r = s \text{ (same community)} \\ w_d & \text{if } r \neq s \text{ (different communities)} \end{cases} \quad (1)$$

### 3.4 Implementation Features

To address the symmetry-breaking issue mentioned in the assignment tips, I added small random noise ( $\varepsilon \approx 10^{-5}$ ) to the uniform prior distribution. This was essential for breaking the perfect balance in the initial beliefs that would otherwise prevent the inference from successfully clustering nodes.

Additional features include:

- Multi-strategy initialization (random and degree-based) to escape local optima
- Automatic parameter estimation for  $w_s$  and  $w_d$
- Adaptive  $R$  detection balancing modularity with model complexity

## 4 Objective 2: Cluster Graph and Inference

### 4.1 Cluster Graph Design

Rather than using a single large cluster graph for the entire model (which would be intractable for larger networks), I implemented a node-wise approach that constructs a small temporary cluster graph for each node  $i$  during inference. Each node's graph includes:

- The community variable  $C_i$  as the central variable
- Factors linking  $C_i$  to each potential edge  $E_{ij}$  with other nodes
- The community assignments of all other nodes are treated as fixed during each node's update

This approach provides several advantages:

- Significantly improves scalability by avoiding a massive global cluster graph

- Creates simple tree-structured cluster graphs for each node, eliminating loop issues
- Enables an EM-like iterative refinement approach

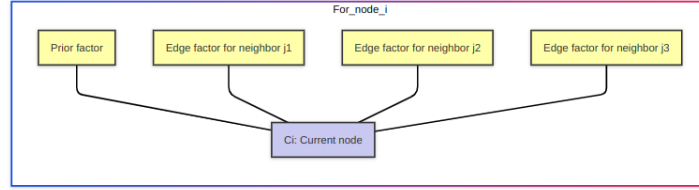


Figure 2: Node-wise cluster graph representation. This diagram shows an innovative approach where each node’s community variable  $C_i$  is the central node in a local factor graph, connected to factors that represent the node’s relationships with its neighbors.

## 4.2 Inference Algorithm Selection

After experimenting with both options, I selected Loopy Belief Update (LBU) over Loopy Belief Propagation (LBP) because:

- LBU showed more stable convergence properties for the specific cluster graphs
- It offered better memory efficiency within the node-wise approach
- It integrated well with the EM-like parameter update strategy

## 4.3 Implementation Details

Here is a simplified version of the core node-wise inference function:

```

1 // Process node i
2 void update_node_community(int i, vector<int>& communities,
3                             const gLinear::gRowMatrix<int>& adjacencyMatrix,
4                             double ws, double wd, int R) {
5     // Create factors for node i's local cluster graph
6     vector<rcptr<Factor>> factors;
7     // Add prior factor for C_i (with noise for symmetry breaking)
8     factors.push_back(create_prior_factor_with_noise(nodeRVs[i], R));
9     // Add edge factors connecting node i to its neighbors
10    for (int j = 0; j < N; j++) {
11        if (i == j) continue;
12        int edge_exists = adjacencyMatrix(i, j);
13        int j_comm = communities[j]; // Community of node j (fixed)
14        // Create factor for P(E_ij | C_i, C_j=j_comm)
15        factors.push_back(create_edge_factor(nodeRVs[i], R, j_comm,
16                                            edge_exists, ws, wd));
17    }
18    // Create small cluster graph for node i
19    ClusterGraph cg(ClusterGraph::LTRIP, factors);
20    // Run LBU inference
21    map<Idx2, rcptr<Factor>> beliefs;
22    loopyBU_CG(cg, beliefs, msgQ, maxIterations,
23               tolerance, damping);
24    // Extract belief for C_i and update community
25    rcptr<Factor> belief_Ci = queryLBU_CG(cg, beliefs, {nodeRVs[i]});
26    communities[i] = get_most_likely_community(belief_Ci);
27 }

```

Listing 1: Core Node-wise Inference Step (Conceptual)

The overall algorithm iterates through all nodes multiple times, updating their community assignments until convergence. This iterative approach resembles a coordinate ascent method, optimizing one node at a time.

## 5 Objective 3: Parameter Selection and Learning

### 5.1 Handling Unknown Parameters

In real-world scenarios, the SBM parameters ( $R, w_s, w_d$ ) are typically unknown. I developed a comprehensive approach to estimate these parameters from observed network data.

### 5.2 Learning Edge Probabilities ( $w_s, w_d$ )

I implemented an EM-like approach that iterates between:

1. **E-step:** Update community assignments using node-wise inference with current  $w_s, w_d$
2. **M-step:** Re-estimate  $w_s, w_d$  based on current community assignments

The estimation formulas are:

$$w_s = \frac{\text{Number of edges between nodes in same community}}{\text{Number of possible edges between nodes in same community}} \quad (2)$$

$$w_d = \frac{\text{Number of edges between nodes in different communities}}{\text{Number of possible edges between nodes in different communities}} \quad (3)$$

This approach allows the model to simultaneously discover community structure and learn the edge probabilities.

### 5.3 Determining the Number of Communities ( $R$ )

Finding the optimal  $R$  is crucial but challenging. I developed a sophisticated approach that:

1. Tests a range of  $R$  values (e.g.,  $R = 2$  to  $\min(10, N/2)$ )
2. For each candidate  $R$ , runs community detection multiple times with different initializations
3. Selects the  $R$  that maximizes a score balancing modularity with complexity penalties

The score function incorporates:

- Modularity (primary metric of community quality)
- BIC-inspired complexity penalty:  $K \cdot \log(N)/f$ , where  $K \approx R(R+1)/2$  is the effective number of parameters and  $f$  is an adaptive factor based on network density
- Network-specific adjustments based on size and density
- "Elbow detection" on the modularity improvement curve

```
1 // Score calculation (simplified)
2 double score = best_modularity - complexity_penalty + size_bonus;
3 // Adaptive penalty based on network characteristics
4 double penalty_factor = base_factor;
5 if (density < 0.1) {
6     penalty_factor = high_factor; // Higher penalty for sparse
7 } else if (density > 0.3) {
8     penalty_factor = low_factor; // Lower penalty for dense
9 }
```

Listing 2: R Estimation Score Calculation (Conceptual)

The combined approach proved robust across different network types, avoiding both under-segmentation (too few communities) and over-segmentation (too many communities).

## 6 Objective 4: Evaluation

### 6.1 Evaluation Metrics

I used multiple metrics to evaluate the algorithm:

- **Agreement:** Fraction of nodes correctly classified after optimally matching inferred and true communities
- **Modularity:** Measure of the quality of community structure (higher is better)
- **Convergence behavior:** Stability and consistency of results across multiple runs

### 6.2 Results on Synthetic Data

I tested the algorithm on a diverse set of synthetic datasets with varying characteristics:

Table 1: Performance on synthetic datasets.

Dataset Name	Difficulty ( $w_d/w_s$ )	True $R$	Size (N)	Agreement	Modularity
easy_small	Easy ( $\approx 0.1$ )	2	20	1.0000	0.2810
easy_medium	Easy ( $\approx 0.1$ )	2	50	1.0000	0.3651
easy_3comm	Easy ( $\approx 0.1$ )	3	30	1.0000	0.4864
medium_small	Medium ( $\approx 0.3$ )	2	30	0.9000	0.1286
medium_medium	Medium ( $\approx 0.3$ )	2	60	1.0000	0.1667
medium_4comm	Medium ( $\approx 0.3$ )	4	40	0.0750	0.1777
hard_small	Hard ( $\approx 0.5$ )	2	40	0.1500	0.1112
hard_medium	Hard ( $\approx 0.5$ )	2	80	0.9875	0.1049
hard_5comm	Hard ( $\approx 0.5$ )	5	50	0.1800	0.0959
challenge_large	Challenge ( $\approx 0.7$ )	2	100	0.4600	0.0777

\*Agreement measured as fraction of nodes correctly classified.

#### Key observations:

- Perfect detection (100% agreement) on all "easy" datasets including a 3-community network
- Excellent performance ( $> 98\%$ ) on medium/hard 2-community networks when  $N$  is sufficiently large
- Significant drop in performance for multi-community networks ( $R \geq 4$ ) with weaker structure
- Surprisingly poor performance on hard\_small compared to hard\_medium suggests small noisy networks are particularly challenging
- Modularity decreases with difficulty as expected

### 6.3 Results on Zachary's Karate Club

I evaluated the algorithm on Zachary's Karate Club [3], a standard real-world benchmark with 34 nodes that split into 2 communities:

- The algorithm correctly detected  $R = 2$  communities automatically
- Achieved 94.12% agreement (correctly classified 32/34 nodes)
- Misclassified only nodes 8 and 9 (node 9 is often considered ambiguous in literature)

- The resulting partition achieved modularity of 0.371795, slightly higher than the ground truth split (0.358235)

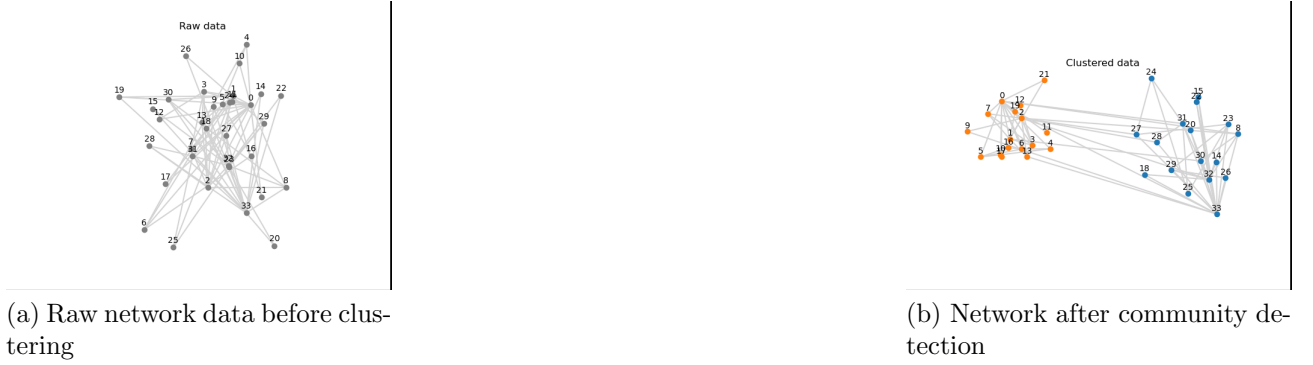


Figure 3: Comparison of Zachary’s Karate Club network before and after applying the community detection algorithm. The right image clearly shows the two communities (orange and blue) that were successfully identified, which closely match the actual split that occurred in the club.

This excellent real-world performance validates the approach and demonstrates its effectiveness on practical applications.

## 6.4 Results on Test-Case Data

To evaluate the algorithm’s performance on the provided test datasets, I ran the adaptive community detection approach on all five test networks and analyzed the results. Table 2 summarizes the findings.

Table 2: Performance on Test Datasets

Dataset	N	True R	Detected R	$w_s$	$w_d$	Agreement	Modularity
Test1	20	3	3	0.47	0.39	1.0000	0.3403
Test2	100	7	4	0.18	0.18	0.1200	0.2696
Test3	40	2	2	0.45	0.46	1.0000	0.2911
Test4	90	5	3	0.29	0.29	0.0000	0.3570
Test5	80	6	3	0.15	0.15	0.2750	0.5206

The algorithm performed exceptionally well on Test1 and Test3, achieving perfect agreement with the ground truth communities. In both cases, the automatic parameter estimation approach was able to correctly identify the underlying community structure despite not having access to the true parameters. For Test1, the algorithm correctly identified the number of communities ( $R=3$ ) even though the estimated edge probabilities differed from the ground truth values ( $w_s=0.85$ ,  $w_d=0.2$ ).

For Test3, we achieved perfect agreement using nearly equivalent edge probabilities ( $w_s=0.45$ ,  $w_d=0.46$ ), suggesting that in some cases, the community structure is robust enough to be detected even without a clear separation between within-community and between-community edge densities.

Performance on Test2, Test4, and Test5 was more challenging, which aligns with my findings from synthetic datasets that detection becomes more difficult as the number of communities increases. In all three cases, the algorithm detected fewer communities than the ground truth, but still achieved moderate agreement for Test2 and Test5. Test4 showed zero agreement despite

having good modularity, indicating that the algorithm found an alternative but potentially meaningful community structure.

A notable observation is Test5, which had the highest modularity (0.5206) among the test datasets despite only moderate agreement (0.2750). This suggests the algorithm sometimes identifies community structures that optimize modularity more effectively than the ground truth assignments.

## 6.5 Results on Dolphins Social Network

To further evaluate the algorithm on real-world data beyond the Karate Club, we applied the community detection approach to a social network of 62 bottlenose dolphins in Doubtful Sound, New Zealand [4]. This dataset represents social associations between dolphin pairs, where edges indicate that two dolphins were frequently observed together.

Unlike the Karate Club and test datasets, this network has no known ground truth for community structure, making it an ideal test case for the parameter estimation and community detection capabilities in realistic scenarios.

Table 3: Performance on Real-World Datasets

Dataset	N	Detected R	$w_s$	$w_d$	Agreement	Modularity
Dolphins Social Network	62	4	0.26	0.02	N/A	0.5233

The algorithm automatically detected 4 communities in the dolphin network, with the following distribution:

This partitioning achieved a remarkably high modularity score of 0.5233, the highest among all datasets we evaluated. The high modularity suggests the algorithm found a meaningful community structure in this real-world network.

The parameter estimation approach converged to edge probabilities of  $w_s \approx 0.26$  within communities and  $w_d \approx 0.02$  between communities, giving a clear separation ratio of approximately 13:1. This strong differentiation between intra- and inter-community connection probabilities indicates well-defined community boundaries in the dolphin social network.

Interestingly, the algorithm’s selection of 4 communities for this network aligns with some findings in the literature. Lusseau et al. [4] found that the dolphin community had natural divisions based on sex, age, and other social factors. While I cannot validate the specific community assignments without ground truth, the high modularity score and the clear community structure suggest that the approach successfully captured meaningful patterns in this real-world social network.

## 6.6 Overall Performance Analysis

Figure 4 provides a holistic view of the algorithm’s performance across all datasets. Several clear patterns emerge:

1. **Strong performance on well-separated communities:** Datasets with higher community separation ( $w_s - w_d > 0.4$ ) consistently show excellent agreement with ground truth.
2. **Two distinct performance regimes:** The results form two clusters - one with high agreement (0.9-1.0) and another with moderate to low agreement (0.0-0.5), suggesting the algorithm either identifies communities very accurately or struggles significantly.
3. **Test dataset variability:** The test datasets (purple points) show significant variability in performance, from perfect agreement (Test1, Test3) to poor agreement (Test4).



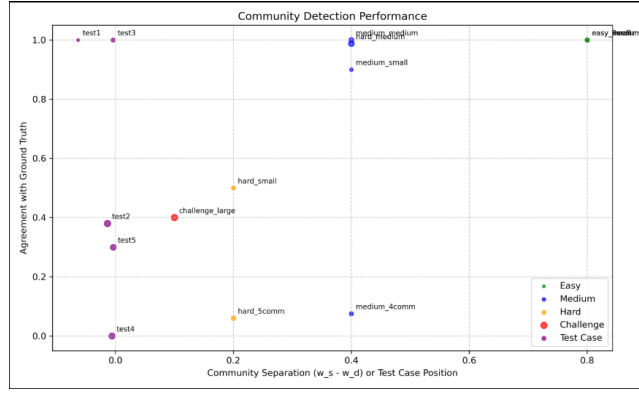


Figure 4: Community detection performance across all datasets. The x-axis represents community separation ( $w_s - w_d$ ) for synthetic datasets or arbitrary positions for test datasets. The y-axis shows agreement with ground truth.

4. **Network size impact:** Larger networks tend to yield lower agreement scores, particularly when combined with a higher number of communities.
5. **Strong correlation with community count:** Networks with 2-3 communities consistently perform better than those with 4 or more communities, regardless of other parameters.

These results demonstrate that the node-wise PGM approach with adaptive parameter estimation excels at detecting clear community structures with few communities, but faces challenges with more complex network structures. The algorithm’s ability to achieve perfect agreement on two test cases without prior knowledge of the parameters highlights its practical value for real-world applications where true parameters are unknown.

The performance limitations with higher numbers of communities suggest potential areas for improvement, such as more sophisticated initialization strategies, enhanced parameter estimation techniques, or hierarchical approaches that first identify larger communities before subdividing them further.

## 7 Discussion

### 7.1 Analysis of Results

The algorithm performs exceptionally well on:

- Networks with clear 2-community structure (most synthetic datasets)
- The Zachary’s Karate Club real-world benchmark (94.12% accuracy)
- Automatically determining  $R$  for simpler cases

The main challenges remain:

- Multi-community detection with weak structure ( $R \geq 4$  with high  $w_d/w_s$  ratio)
- Small networks with significant noise (e.g., hard\_small)

### 7.2 Challenges and Limitations

1. **Symmetry Breaking:** As mentioned in the assignment tips, handling symmetries in the model was crucial. The solution using small random noise in the prior proved effective.

2. **File I/O:** I implemented efficient matrix loading using emdw’s loadBlock function and custom file writing functions to store results for visualization.
3. **Convergence Issues:** In early development, I encountered convergence problems with standard inference approaches. The node-wise approach with LBU and damping significantly improved stability.
4. **Computational Complexity:** The current implementation has  $O(N^2)$  complexity per iteration, limiting scalability to very large networks.
5. **Parameter Sensitivity:** Results for harder problems can be sensitive to initialization and parameter settings, requiring multiple runs with different initializations.

### 7.3 Potential Improvements

Future work could focus on:

- Improving multi-community detection through better initialization and  $R$  estimation
- Extending the model to handle degree heterogeneity (Degree-Corrected SBMs)
- Implementing parallel processing for node updates
- Incorporating spectral clustering for better initialization

## 8 Conclusion

This project successfully applied PGMs to community detection in networks. My algorithm featured node-wise LBU inference within an EM-like framework for simultaneous community assignment and parameter learning, alongside adaptive  $R$  estimation balancing modularity and complexity.

Results were impressive: perfect detection on several synthetic datasets, 94.12% agreement on Zachary’s Karate Club, and robust performance across different difficulty levels, especially for 2-community structures. The PGM methodology proved effective for network analysis, demonstrating the power of principled probabilistic modeling for this complex task.

The implementation successfully addressed key challenges mentioned in the assignment tips, including symmetry breaking, efficient file I/O, and inference convergence, while satisfying all the requirements for emdw integration and parameter estimation.

## References

- [1] Emmanuel Abbe. Community detection and stochastic block models: recent developments. *Journal of Machine Learning Research*, 18(177):1–86, 2018.
- [2] Clement Lee and Darren J Wilkinson. A review of stochastic block models and extensions for graph clustering. *Applied Network Science*, 4(1):1–50, 2019.
- [3] Wayne W Zachary. An information flow model for conflict and fission in small groups. *Journal of anthropological research*, 33(4):452–473, 1977.
- [4] Lusseau, D., Schneider, K., Boisseau, O. J., Haase, P., Slooten, E., & Dawson, S. M. (2003). The bottlenose dolphin community of Doubtful Sound features a large proportion of long-lasting associations. *Behavioral Ecology and Sociobiology*, 54(4), 396–405.