

Selecting the Number of Communities for Weighted Degree-Corrected Stochastic Block Models

Yucheng Liu Xiaodong Li

University of California, Davis



Introduction

- Weighted networks often reveal more refined community structure than corresponding unweighted ones.
- We want to study how to consistently estimate the number of communities in weighted networks under mild conditions.

Contributions:

- 1 We propose the weighted DCSBM, which is a generic model for weighted networks without modeling the likelihood.
- 2 We propose a stepwise testing procedure in selecting the number of communities under the weighted DCSBM and prove the consistency of the procedure under mild conditions.
- 3 We generalize the Nonsplitting Property of SCORE to weighted DCSBM.
- 4 Simulations on both synthetic and real-world weighted networks show empirical consistency of our proposed procedure.

Weighted DCSBM

- **Key difference from standard DCSBM:** Instead of specifying the exact likelihood, the model only specifies the first two moments.
- **Parameters:** θ_i is the heterogeneity parameter; $\boldsymbol{\pi}_i \in \mathbb{R}^K$ indicates the community belonging of node i ; \mathbf{B} is the $K \times K$ symmetric community connectivity matrix.
- **First moment:** $\mathbb{E}[A_{ij}] := M_{ij} = \theta_i \theta_j \boldsymbol{\pi}_i^\top \mathbf{B} \boldsymbol{\pi}_j$.
- **Second moment:** $V_{ij} := \text{var}(A_{ij}) = \nu(M_{ij})$, where ν is known from the underlying distribution.

Estimation of parameters:

- Node community belonging is estimated by SCORE.
- $\hat{\theta}_i^{(m)} := \frac{\sqrt{(\hat{\mathbf{1}}_k^{(m)})^\top \mathbf{A} \hat{\mathbf{1}}_k^{(m)}}}{(\hat{\mathbf{1}}_k^{(m)})^\top \mathbf{A} \mathbf{1}_n} d_i$, where d_i is the degree of node i .
- $\hat{B}_{kl}^{(m)} := \frac{(\hat{\mathbf{1}}_k^{(m)})^\top \mathbf{A} \hat{\mathbf{1}}_l^{(m)}}{\sqrt{(\hat{\mathbf{1}}_k^{(m)})^\top \mathbf{A} \hat{\mathbf{1}}_k^{(m)}} \sqrt{(\hat{\mathbf{1}}_l^{(m)})^\top \mathbf{A} \hat{\mathbf{1}}_l^{(m)}}}$.

Assumption 1

Denote $\theta_{\max} = \max\{\theta_1, \dots, \theta_n\}$ and $\theta_{\min} = \min\{\theta_1, \dots, \theta_n\}$. c_0 is a small constant. We assume the following conditions hold:

- [Fixed rank] The true number of communities K is fixed.
- [Balancedness]
$$\min_{1 \leq k \leq K} \frac{n_k}{n} \geq c_0 \quad \text{and} \quad \frac{\theta_{\min}}{\theta_{\max}} \geq c_0.$$
- [Sparseness]
$$\frac{1}{c_0} \geq \theta_{\max} \geq \theta_{\min} \geq \frac{\log^3 n}{\sqrt{n}}.$$
- [Community connectivity] The $K \times K$ matrix \mathbf{B} is fixed, and its entries and eigenvalues satisfy
$$\begin{cases} B_{kk} = 1 & \text{for } k = 1, \dots, K, \\ c_0 \leq B_{kl} \leq 1 & \text{for } 1 \leq k, l \leq K, \\ \lambda_1(\mathbf{B}) > |\lambda_2(\mathbf{B})| \geq \dots \geq |\lambda_K(\mathbf{B})| \geq c_0 > 0. \end{cases}$$
- [Variance-mean function] The function $\nu(\cdot)$ satisfies $c_0 \mu \leq \nu(\mu) \leq \mu/c_0$ and $\nu(\cdot)$ is $1/c_0$ -Lipschitz.
- [Bernstein condition] For any $i \leq j$ and any integer $p \geq 2$, there holds

$$\mathbb{E}[|A_{ij} - M_{ij}|^p] \leq \left(\frac{p!}{2}\right) R(c_0)^{p-2} \nu(M_{ij}),$$

where $R(c_0)$ is a constant only depending on c_0 .

Our Algorithm

SVPS: Stepwise Variance Profile Scaling

For $m = 1, 2, \dots$:

- 1 Group nodes into m distinct communities using SCORE.
- 2 Obtain estimated mean adjacency matrix $\widehat{\mathbf{M}}^{(m)}$ by fitting DCSBM and derive the estimated variance profile matrix $\widehat{\mathbf{V}}^{(m)}$ using the variance-mean relationship.
- 3 Find the scaling matrix $\widehat{\boldsymbol{\Psi}}^{(m)}$ such that $\widehat{\boldsymbol{\Psi}}^{(m)} \widehat{\mathbf{V}}^{(m)} \widehat{\boldsymbol{\Psi}}^{(m)}$ is doubly stochastic (every row sum equals 1).
- 4 Obtain test statistic $T_{n,m} = \left| \lambda_{m+1} \left(\left(\widehat{\boldsymbol{\Psi}}^{(m)} \right)^{\frac{1}{2}} \mathbf{A} \left(\widehat{\boldsymbol{\Psi}}^{(m)} \right)^{\frac{1}{2}} \right) \right|$.
Stop the procedure if $T_{n,m} < 2 + \epsilon$ and obtain $\hat{K} = m$.

Theoretical Results

Theorem (Null Case). If we implement SVPS with $m = K$ and SCORE for spectral clustering, then for any fixed $c_0 > 0$ in Assumption 1, as $n \rightarrow \infty$, we have $T_{n,m} \leq 2 + o_P(1)$.

Theorem (Underfitting Case). If we implement SVPS with $m < K$ and SCORE for spectral clustering, then for any fixed $c_0 > 0$ in Assumption 1, as $n \rightarrow \infty$, we have $T_{n,m} \xrightarrow{P} \infty$.

Definition (Nonsplitting Property (Jin et al., 2022)). The estimated communities of a network satisfy the Nonsplitting Property if the true communities are a refinement of the estimated ones.

Lemma (Nonsplitting Property). Under Assumption 1, with any fixed $c_0 > 0$, for any fixed $m \leq K$, SCORE satisfies the NSP with probability $1 - O(n^{-3})$.

Experiments

Generating mechanism of synthetic networks:

- Underlying generating distributions: Poisson, binomial and negative binomial.
- $B_{kl} = \rho (1 + r \times \mathbf{1}_{\{k=l\}})$.

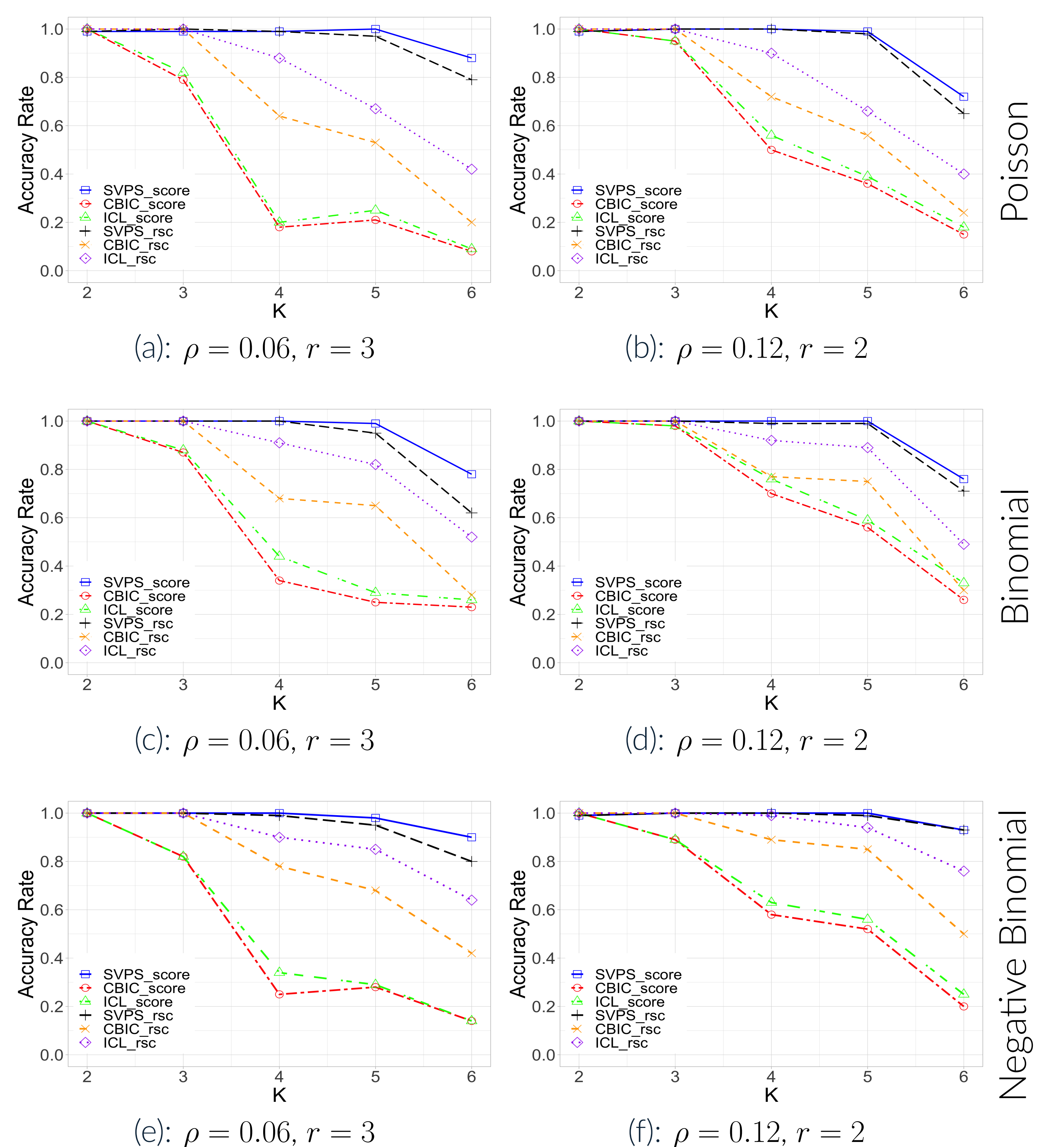


Figure 1: Accuracy rate of SVPS and other methods for comparison.