

# Order Estimation in Locally Independent Sub-Exponential Mixtures by Stepwise Goodness-of-Fit

Wenjin Zhang<sup>†</sup>

University of Science and Technology of China

## Abstract

Estimating the number of components in finite mixture models (FMMs) is a critical problem in statistical methodology. While traditional methods often focus on Gaussian or sub-Gaussian mixtures, real-world data usually exhibit heavier tails. This paper extends the Stepwise Goodness-of-Fit (StGoF) method (Jin et al., 2022) to locally independent sub-exponential mixtures. Our method has two main advantages: it avoids the need for precise parameter estimation, enabling faster computation, and it supports mixed-type data. We establish theoretical guarantees for exact recovery and asymptotic consistency, and demonstrate the efficiency and robustness of our approach through simulations and real-world applications. This work introduces a new method for estimating the order of locally-independent sub-exponential mixtures, paving the way for further research.

**Keywords:** Mixture Models, Sub-Exponential Distributions, Stepwise Goodness-of-Fit, Clustering.

## 1 Introduction

Finite mixture models (FMMs) have been extensively studied in the statistical literature, forming a cornerstone of modern statistical methodology (McLachlan and Peel, 2004; McLachlan et al., 2019; Bouguila and Fan, 2020). Representing complex distributions as weighted sums of simpler components, FMMs provide a flexible and robust framework for modeling heterogeneous data. This flexibility makes FMMs particularly valuable in uncovering latent structures, clustering observations, and performing density estimation. Their applications span numerous fields, including machine learning (Goodfellow et al., 2016), genetics (Bechtel et al., 1993), and medical research (Schlattmann, 2009). With a rich history of development, FMMs demonstrate both theoretical depth and practical versatility across disciplines.

Building upon the versatility and widespread applications of finite mixture models (FMMs), a critical aspect of their practical implementation is determining the correct number of components, or the model’s order. Accurate estimation of the order is essential to ensure model interpretability and estimation efficiency. An underestimated model fails to capture the complexity of the data, while an over-specified model introduces unnecessary complexity, deteriorating estimation rates and parameter reliability. These challenges have spurred extensive research into methods for estimating the order of FMMs.

---

<sup>†</sup>Research conducted under the guidance of Professor Yuqi Gu at Columbia University. This is a working manuscript that is still being updated. Final authorship: Wenjin Zhang and Yuqi Gu.

Numerous methods have been proposed for estimating the order of FMMs. Likelihood-based approaches, such as hypothesis testing (McLachlan, 1987; Dacunha-Castelle and Gassiat, 1999; Liu and Shao, 2003) and the EM-test (Chen and Li, 2009; Li and Chen, 2010), focus on evaluating nested models and typically assume prior knowledge of a candidate order. Information criteria, including AIC (Akaike, 1974) and BIC (Schwarz, 1978), are among the most widely used techniques, with BIC being particularly favored for estimating the number of mixture components (Leroux, 1992; Keribin, 2000; McLachlan and Peel, 2004). Extensions, such as the Integrated Completed Likelihood (Biernacki et al., 2000) and Singular BIC (Drton and Plummer, 2013), have been proposed to address the challenges posed by non-regular models. More recently, methods such as Group-Sort-Fuse (Manole and Khalili, 2021) and Evidence Lower Bound maximization (Wang and Yang, 2024) have further advanced this field. Despite their popularity, likelihood-based methods typically require iterative algorithms like the Expectation-Maximization (EM) algorithm to estimate parameters, which can be computationally expensive and slow, particularly in high-dimensional settings. Alternatively, minimum-distance-based methods (Chen and Kalbfleisch, 1996; James et al., 2001; Woo and Sriram, 2006; Heinrich and Kahn, 2018) minimize discrepancies between observed data and candidate models, offering a flexible alternative to likelihood-based techniques. However, these methods also depend on pre-specified parametric forms for the component distributions, limiting their ability to address scenarios where the data's underlying distribution is unknown or mixed.

In addition to the general methods discussed above, much of the existing literature has focused on specific types of mixture models, particularly Gaussian mixture models (GMMs) and sub-Gaussian mixture models. GMMs have been extensively studied for their mathematical tractability and wide applicability in clustering and parameter estimation under separation conditions (Vempala and Wang, 2004; Ndaoud, 2018; Zhang and Zhou, 2021; Chen and Yang, 2021). sub-Gaussian mixture models, on the other hand, address settings with lighter-tailed distributions, providing strong theoretical guarantees for clustering and recovery in high-dimensional scenarios (Mixon et al., 2017; Srivastava et al., 2019; Cai and Zhang, 2018; Abbe et al., 2022).

However, in fields such as finance and economics, data often exhibit heavier tails that cannot be adequately captured by Gaussian or sub-Gaussian models. sub-exponential mixture models, with their ability to accommodate such heavy-tailed distributions, offer a more suitable framework for these applications. Despite their potential, research on sub-exponential mixture models remains limited (Dreveton et al., 2024), highlighting the need for further exploration in both theoretical development and practical methodology.

In a remarkable paper, Jin et al. (2022) propose a stepwise Goodness-of-Fit (StGoF) method to estimate the number of communities in degree-corrected block models (DCBM). Their work introduces a stepwise algorithm that alternates between a community detection step and a Goodness-of-Fit step for  $m = 1, 2, \dots$ . The core idea of this framework is highly adaptable and can be applied to the context of interest in this paper. Specifically, we extend the StGoF framework to estimate the order of mixtures, where the variables are assumed to be conditionally independent given their membership labels, and the noise follows a distribution with sub-exponential tails. While the assumption of conditional (or local) independence might seem restrictive at first glance, it actually accommodates a broad range of well-known models, such as Spherical Gaussian Mixture Models and Latent Class Models. As such, this assumption not only preserves the generality of our approach but also ensures analytical tractability.

Our modifications retain the core structure of the algorithm, alternating between clustering and goodness-of-fit steps, but we adapt both steps to account for the properties of sub-exponential noise. At each iteration, the clustering step applies a clustering algorithm to the data, identifying the current partitioning of observations into clusters. This is followed by a GoF step, where we calculate a test score based on the clustering results from the previous step. While this approach

is presented in the context of sub-exponential mixture models for clarity, it is worth noting that the framework is general enough to accommodate broader applications beyond locally independent sub-exponential mixtures, such as Weibull Distribution with  $k \in (0, 1)$ .

In this study, we extend the StGoF framework, originally developed for network analysis, to the setting of mixture models. Although the core framework remains the same, its implications are far-reaching. Firstly, our findings demonstrate that the refitted quadrilateral test statistics introduced by Jin et al. (2022) are not limited to degree-corrected block models but are also applicable to the mixture models examined herein. Indeed, the scope of this method is broader, offering significant applicability to a variety of problems that involve estimating the number of clusters or components. This issue remains an active area of investigation in our ongoing research.

Secondly, our approach is not contingent upon the specific form of the data distribution. Given an observed  $p$ -dimensional dataset, each dimension can follow an arbitrary distribution—such as Gaussian, Poisson, or Gamma—provided that each component is sub-exponential. This contrasts with most existing methods, which typically require prior knowledge of the distributional assumptions. As previously noted, many likelihood-based techniques rely on explicit assumptions regarding the distribution of each variable. In the context of mixture models, deviations from these parametric assumptions can lead to suboptimal clustering outcomes (Foss et al., 2019). For real-world datasets, such as high-dimensional mixed-type data, it is often challenging to determine the distribution type of each component within a  $p$ -dimensional variable. In the case of the gap statistics method (Tibshirani et al., 2000; Hennig and Lin, 2015), obtaining better results for mixed-type data requires carefully selecting an appropriate distance metric within the internal criterion. Although considerable research has focused on hybrid distance metrics (Gower, 1971; Huang, 1998; Ahmad and Dey, 2011; Hennig and Liao, 2013), the majority of these studies typically consider only numerical and categorical data, without distinguishing between different types of numerical data, such as continuous and count data. In such cases, these criteria often rely on simple distance measures to assess the quality of clustering. Given the heterogeneity of different distribution types, the performance of such methods may not always be optimal, especially in complex datasets with diverse data types. In this regard, our methodology offers a significant advantage by allowing for direct analysis of more complex mixed-type data.

Finally, another advantage of our approach is its computational efficiency. Unlike likelihood-based methods, which require iterative optimization to estimate parameters, our method does not involve such procedures, leading to significantly faster computations. Additionally, the prediction strength method (Tibshirani and Walther, 2005; Dudoit and Fridlyand, 2002; Volkovich et al., 2011) requires either cross-validation or resampling, which similarly incurs substantial computational costs. In contrast, the primary computational costs in StGoF arise from two steps: a single application of k-means clustering and the computation of a matrix statistic at each step. These operations enable much faster result computation compared to methods that require iterative optimization and cross validation.

**Organization.** The remainder of the paper is organized as follows. Section 2 presents the mathematical formulation of our method, including the stepwise goodness-of-fit test and matrix correction procedure. In Section 3, we provide theoretical results on the consistency of our estimator. Section 4 illustrates the performance of our method through simulations and real data applications. Section 5 concludes with a discussion of future research directions. The Supplementary Material contains all technical proofs of the theoretical results.

**Notation.** For a matrix  $\mathbf{A}$ , the notation  $\mathcal{S}(\mathbf{A})$  refers to the symmetric dilation of  $\mathbf{A}$ , which is

defined as:

$$\mathcal{S}(\mathbf{A}) = \begin{pmatrix} 0 & \mathbf{A} \\ \mathbf{A}^\top & 0 \end{pmatrix}.$$

The  $i$ -th row of  $\mathbf{A}$  is denoted by  $r_i(\mathbf{A})$ , and the  $i$ -th largest singular value of  $A$  is represented by  $\sigma_i(\mathbf{A})$ . The notation  $\mathbf{A}_{1:m}$  refers to the submatrix consisting of the first  $m$  columns of  $\mathbf{A}$ . We denote the operator norm of  $\mathbf{A}$  by  $\|\mathbf{A}\|$ , and the infinity norm  $\|\mathbf{A}\|_\infty$  as the maximum absolute value of any element in  $\mathbf{A}$ . When applied to the noise matrix  $E$ , the symmetric dilation is denoted by  $\mathbf{W} = \mathcal{S}(E)$ . For the signal matrix  $\mathbf{P}$ , the singular values are ordered as  $\sigma_1 \geq \dots \geq \sigma_K$ .

## 2 Model and Methodology

In this paper, we study a mixture model consisting of  $K$  clusters, denoted  $C_1, C_2, \dots, C_K$ , with each cluster centered at  $\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \dots, \boldsymbol{\theta}_K^* \in \mathbb{R}^{1 \times p}$ . The minimum separation between any two distinct cluster centers is defined as  $\Delta = \min_{a \neq b} \|\boldsymbol{\theta}_a^* - \boldsymbol{\theta}_b^*\|$ . The assignment of the  $n$  observations to these clusters is described by a cluster assignment vector  $\mathbf{z} \in \{1, \dots, K\}^n$ , where  $z_i$  indicates the cluster to which observation  $\mathbf{X}_i \in \mathbb{R}^{1 \times p}$  belongs. Let  $\mathbf{Z}$  represent the matrix  $(z_1^\top \ \dots \ z_n^\top)^\top$ . For convenience, we define  $N = n + p$ .

Each observation is generated according to the following model:

$$\mathbf{X}_i = \boldsymbol{\theta}_{z_i}^* + \boldsymbol{\epsilon}_i, \quad (1)$$

where  $\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_n \in \mathbb{R}^{1 \times p}$  are noise terms. The observations  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are stacked row-wise into a data matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , which can be expressed in matrix form as:

$$\mathbf{X} = \mathbf{P} + \mathbf{E},$$

where  $\mathbf{P} := (\boldsymbol{\theta}_{z_1}^*; \boldsymbol{\theta}_{z_2}^*; \dots; \boldsymbol{\theta}_{z_n}^*)$  is the signal matrix containing the true cluster centers, and  $\mathbf{E} := (\boldsymbol{\epsilon}_1; \dots; \boldsymbol{\epsilon}_n)$  is the noise matrix. We can rewrite  $\mathbf{P}$  as  $\mathbf{Z}\boldsymbol{\Theta}^\top$ , where  $\boldsymbol{\Theta}$  denotes the matrix  $((\boldsymbol{\theta}_1^*)^\top \ \dots \ (\boldsymbol{\theta}_K^*)^\top)^\top$ .

**Remark 1.** *The model assumptions in this paper largely follow those proposed in Zhang and Zhou (2022). Notably, our assumptions differ from the standard ones commonly adopted in finite mixture models (FMMs). Typically, for FMMs, the assumptions are formulated as follows: Denoting  $\mathbf{z}^* \in [k]^n$  as the vector of cluster assignments, a mixture model assumes that the  $n$  observed data points  $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathcal{X}^n$ , where  $\mathcal{X} \subset \mathbb{R}^d$ , are independently generated such that*

$$\forall i \in [n] : \mathbf{X}_i \mid \mathbf{z}_i^* \sim f_{\mathbf{z}_i^*},$$

where  $f_1, \dots, f_k$  are  $k$  probability distributions over  $\mathcal{X}$  (see Dreveton et al. (2024) for a detailed discussion).

In contrast, our model is more general given membership labels. Specifically, we do not require the noise distributions within the same cluster to be identical. Furthermore, our model explicitly defines the cluster centers  $\boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_K^*$ . This explicit specification is crucial because the success of our clustering method depends on exact recovery of the clusters. Consequently, information about the minimum distance between clusters,  $\Delta$ , is indispensable for achieving exact recovery. This requirement necessitates the structural design of our model, which incorporates explicit centers and a clear delineation of cluster separations.

Before presenting the proposed algorithm, we introduce additional assumptions about the noise terms, which will be used throughout the remainder of this paper.

**Assumption 1.** The components of the noise vectors  $\epsilon_1, \dots, \epsilon_n$ , specifically the elements  $\epsilon_{i,j}$  (where  $1 \leq i \leq n$  and  $1 \leq j \leq p$ ), are assumed to be mutually independent. Each  $\epsilon_{i,j}$  follows a sub-exponential distribution satisfying  $\max_{i,j} \text{Var}(\epsilon_{i,j}^2) \leq \max_{i,j} \|\epsilon_{i,j}\|_{\psi_1}^2 = \sigma^2$  and  $\text{Var}(\epsilon_{i,j}) \geq \tau^2$ , where  $\sigma$  and  $\tau$  are fixed constants.

**Remark 2.** At first glance, the assumption  $\text{Var}(\epsilon_{i,j}) \geq \tau^2$  may appear somewhat unconventional. In fact, this condition is used solely in the proof of Lemma S.1 to establish the asymptotic normality of a specific statistic. However, it is worth noting that Lemma S.1 is not essential for proving our main result, Theorem 2, and thus this assumption could be omitted without affecting the validity of the main theorem. We retain this assumption here because Lemma S.1 is a stronger result that could offer additional insights and potentially support further developments in this method.

The procedure in our case is summarized in Algorithm 1. Specifically, in the clustering step of Algorithm 1, we employ a standard spectral clustering method. The statistic  $z_\alpha$  is the  $\alpha$ -quantile of the standard normal distribution  $N(0, 1)$ .

---

**Algorithm 1:** Stepwise Goodness-of-Fit

---

**Require:** Data matrix  $\mathbf{X}$ ,  $\sigma$ ,  $\beta$  (initialize  $m = 1$ )

**Ensure:** Estimated number of components  $\hat{K}_\alpha$

- 1: **Clustering:** Perform top- $m$  SVD on  $\mathbf{X}$ ; apply  $k$ -means clustering with  $m$  clusters to rows of  $(\mathbf{U}_{\mathbf{X}})_{1:m} \in \mathbb{R}^{n \times m}$ ; obtain membership matrix  $\widehat{\mathbf{Z}}^{(m)}$ .
- 2: **Matrix Correction:** Estimate the signal matrix  $\widehat{\mathbf{P}}^{(m)} = \widehat{\mathbf{Z}}^{(m)} \left( (\widehat{\mathbf{Z}}^{(m)})^\top \widehat{\mathbf{Z}}^{(m)} \right)^{-1} (\widehat{\mathbf{Z}}^{(m)})^\top \mathbf{X}$ , and apply symmetric dilation to  $\widehat{\mathbf{P}}^{(m)}$ , denoted as  $\mathcal{S}(\widehat{\mathbf{P}}^{(m)})$ .
- 3: **Goodness-of-Fit Test:** Compute

$$Q_N^{(m)} = \sum_{\substack{i_1, i_2, i_3, i_4 \\ \text{distinct}}} \left( \mathcal{S}(\mathbf{X})_{i_1, i_2} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_1, i_2} \right) \left( \mathcal{S}(\mathbf{X})_{i_2, i_3} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_2, i_3} \right) \\ \cdot \left( \mathcal{S}(\mathbf{X})_{i_3, i_4} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_3, i_4} \right) \left( \mathcal{S}(\mathbf{X})_{i_4, i_1} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_4, i_1} \right).$$

Calculate  $C_N = 2\sigma^8 n^{\frac{\beta}{2}} p^{\frac{\beta}{2}}$ , and compute the test score  $\phi_N^{(m)} = \frac{Q_N^{(m)}}{\sqrt{C_N}}$ .

- 4: **Termination:** If  $\phi_N^{(m)} > z_\alpha$ , increment  $m$  and repeat Steps 1–3; otherwise, terminate and set  $\hat{K}_\alpha = m$ .
  - 5: **return**  $\hat{K}_\alpha$
- 

As we will demonstrate in Theorem 1 in Section 3, this algorithm achieves exact recovery provided the following conditions are satisfied:

**Assumption 2.** There exists a constant  $\alpha_0 \in (0, 1)$  such that  $|C_k| \geq \alpha_0 n$  for  $\forall 1 \leq k \leq K$ .

**Assumption 3.**  $\sigma_K = \omega(\sigma(\sqrt{n} + \sqrt{p}))$ ,  $\Delta = \omega(\sigma\sqrt{n})$ .

These assumptions ensure a sufficiently high signal-to-noise ratio and balanced community sizes, which are critical for the algorithm to achieve exact recovery of the true community structure.

Additionally, the GoF step is based on the refitted quadrilateral test statistic proposed in the original paper, but we make slight modifications to adapt it to our case. Specifically, in Jin et al. (2022), the data matrix is square. In contrast, we handle non-square matrices by applying symmetric dilation, after which the test score is computed based on the resulting square matrix.

Furthermore, we make adjustments to the final test score to address the challenges specific to our setting.

Our algorithm has a distinct advantage in that it does not rely on the specific type of noise distribution. Unlike traditional likelihood-based methods, which require explicit likelihood computations, our approach is likelihood-free. This eliminates the need for complex likelihood evaluations, significantly improving computational efficiency. Moreover, traditional methods often depend on identifying an appropriate distribution for the data, as the performance of these methods heavily relies on the correctness of the assumed distribution. In contrast, our algorithm bypasses this requirement, allowing it to perform robustly without explicit distributional assumptions.

Additionally, this flexibility enables our method to handle mixed-type data, where the noise can come from heterogeneous or mixed distributions. This generalization broadens the applicability of our approach to a wider range of practical scenarios, where data types and noise characteristics are not uniform or well-defined.

Compared to the degree-corrected block models (DCBM) analyzed in the original paper ([Jin et al., 2022](#)), our model, while not incorporating degree heterogeneity, features a parameter space defined by a  $K \times p$  matrix. This parameterization is substantially larger than the  $K \times K$  probability matrix in DCBM, introducing significant complexity. As a result, deriving a null distribution that adheres to  $N(0, 1)$  becomes substantially more challenging. Additionally, our noise is sub-exponential with an unknown variance  $\sigma$ , adding another layer of difficulty compared to the simpler Bernoulli-distributed noise assumed in SBM. These challenges complicate the construction of an accurate hypothesis testing framework. Consequently, we adopt the current form of the test score, while recognizing that further improvements remain an interesting direction for future research.

In practice, the variance  $\sigma^2$  is often unknown and cannot be directly determined in many cases, necessitating adjustments to the test statistic to account for this uncertainty. However, we retain the current form of the test score because, for certain types of noise,  $\sigma$  can be explicitly determined. For instance, if the noise follows a bounded distribution within the interval  $[a, b]$ , such as a Bernoulli distribution, the variance  $\sigma^2$  can be expressed as  $\frac{(b-a)^2}{4}$ . Similarly, for a Poisson distribution, the variance is equal to the mean, allowing it to be estimated directly from the data matrix. These examples demonstrate that the method accommodates specific cases where noise characteristics are known, ensuring the validity of the test under such conditions.

In general, estimating  $\sigma$  is required before running this algorithm. However, it is evident that estimating  $\sigma$  becomes impossible if the noise terms are drawn from completely distinct distributions. Therefore, we adopt a common assumption in clustering mixture models, where noise terms associated with the same cluster center come from the same distribution. Thus, Assumption 1 is modified to:

**Assumption 1'.** *The components of the noise vectors  $\epsilon_1, \dots, \epsilon_n$ , specifically the elements  $\epsilon_{i,j}$  (where  $1 \leq i \leq n$  and  $1 \leq j \leq p$ ), are assumed to be mutually independent. For  $j = 1, \dots, p$ , the noise terms  $\epsilon_{1,j}, \dots, \epsilon_{n,j}$  are independently drawn from  $K$  distinct sub-exponential distributions  $F_1, \dots, F_K$ . Each noise term satisfies  $\epsilon_{i,j} \sim F_{z_i}$ , where  $z_i$  represents the cluster assignment for the  $i$ -th observation. The sub-exponential distributions  $F_k$  (for  $k = 1, \dots, K$ ) satisfy  $\sup_{k=1, \dots, K} \|x \sim F_k\|_{\psi_1} \leq \sigma$ , where  $\sigma$  is a fixed constant. Moreover, the variance of the noise terms satisfies  $\text{Var}(\epsilon_{i,j}) \geq \tau^2$  for all  $i, j$ , where  $\tau > 0$  is a constant.*

Under this assumption, we can use any standard variance estimator to estimate  $\sigma$  in each iteration based on the clustering result from step (a). This estimate is then used to update  $C_N$  in the algorithm, replacing it with  $\hat{C}_N = 2\hat{\sigma}^8 n^{\frac{\beta}{2}} p^{\frac{\beta}{2}}$ .

For example, we may estimate  $\hat{\sigma}$  as follows:

$$\hat{\sigma} = \max_{1 \leq j \leq p} \max_{1 \leq i \leq m} \left\{ \frac{\sum_{l \in C_i^{(m)}} (X_{l,j} - \frac{1}{|C_i^{(m)}|} \sum_{s \in C_i^{(m)}} X_{s,j})^2}{|C_i^{(m)}| - 1} \right\},$$

where  $C_1^{(m)}, \dots, C_m^{(m)}$  are pseudo communities obtained by step (a).

Given that the community detection method achieves exact recovery and the community sizes are balanced under Assumption 1, this estimator provides a proper estimate of  $\sigma$  when  $m = K$ .

If  $m \neq K$ , the estimator may deviate from the true upper bound of the sub-exponential noise norms. However, any deviation will only lead to an underestimate of the true  $\sigma$ . As we will show in Theorem 3, this underestimation does not affect the final asymptotics of  $\phi_N^{(m)}$ , and thus consistency will still hold.

Thus, we propose an alternative formulation of the Stepwise Goodness-of-Fit procedure in Algorithm 2, which is applicable in the absence of prior knowledge regarding  $\sigma$ . Furthermore, the estimation of  $\sigma$  can be adapted to any suitable method of estimation.

---

**Algorithm 2:** Modified Stepwise Goodness-of-Fit

---

**Require:** Data matrix  $\mathbf{X}$ ,  $\beta$  (initialize  $m = 1$ )

**Ensure:** Estimated number of components  $\hat{K}_\alpha$

- 1: **Clustering:** Perform top- $m$  SVD on  $\mathbf{X}$ ; apply  $k$ -means clustering with  $m$  clusters to rows of  $(\mathbf{U}_x)_{1:m} \in \mathbb{R}^{n \times m}$ ; obtain membership matrix  $\widehat{\mathbf{Z}}^{(m)}$ .
- 2: **Matrix Correction:** Estimate the signal matrix  $\widehat{\mathbf{P}}^{(m)} = \widehat{\mathbf{Z}}^{(m)} \left( (\widehat{\mathbf{Z}}^{(m)})^\top \widehat{\mathbf{Z}}^{(m)} \right)^{-1} (\widehat{\mathbf{Z}}^{(m)})^\top \mathbf{X}$ , and apply symmetric dilation to  $\widehat{\mathbf{P}}^{(m)}$ , denoted as  $\mathcal{S}(\widehat{\mathbf{P}}^{(m)})$ .
- 3: **Goodness-of-Fit Test:** Compute

$$Q_N^{(m)} = \sum_{\substack{i_1, i_2, i_3, i_4 \\ \text{distinct}}} \left( \mathcal{S}(\mathbf{X})_{i_1, i_2} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_1, i_2} \right) \left( \mathcal{S}(\mathbf{X})_{i_2, i_3} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_2, i_3} \right) \cdot \left( \mathcal{S}(\mathbf{X})_{i_3, i_4} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_3, i_4} \right) \left( \mathcal{S}(\mathbf{X})_{i_4, i_1} - \mathcal{S}(\widehat{\mathbf{P}}^{(m)})_{i_4, i_1} \right).$$

Calculate  $\widehat{C}_N = 2\widehat{\sigma}^8 n^{\frac{\beta}{2}} p^{\frac{\beta}{2}}$ , where  $\widehat{\sigma} = \max_{1 \leq j \leq p} \max_{1 \leq i \leq m} \left\{ \frac{\sum_{l \in C_i^{(m)}} (X_{l,j} - \frac{1}{|C_i^{(m)}|} \sum_{s \in C_i^{(m)}} X_{s,j})^2}{|C_i^{(m)}| - 1} \right\}$ ,

and compute the test score  $\frac{Q_N^{(m)}}{\sqrt{\widehat{C}_N}}$ .

- 4: **Termination:** If  $\phi_N^{(m)} > z_\alpha$ , increment  $m$  and repeat Steps 1–3; otherwise, terminate and set  $\hat{K}_\alpha = m$ .
  - 5: **return**  $\hat{K}_\alpha$
- 

### 3 Theoretical Guarantee

Following a similar approach as in the proof of Jin et al. (2022), we first need to establish the Non-Splitting Property (NSP) of the spectral clustering method. This property not only ensures that our method can achieve exact recovery of communities but also serves as a foundational result for analyzing the clustering behavior in underfitting cases. To proceed, we introduce the following definitions.

**Definition 1.** Fix  $K > 1$  and  $m \leq K$ . We say that a realization of the  $n \times m$  matrix of estimated labels  $\widehat{\mathbf{Z}}^{(m)}$  satisfies the NSP if for any pair of nodes in the same (true) community, the estimated community labels are the same (i.e., each community in  $\mathbf{Z}$  is contained in a community in the realization of  $\widehat{\mathbf{Z}}^{(m)}$ ). When this happens, we write  $\mathbf{Z} \preceq \widehat{\mathbf{Z}}^{(m)}$ .

From the definition, it follows that if a clustering method satisfies the NSP, then when it is tasked with partitioning the data into  $m$  clusters ( $m < K$ ), the resulting clusters are formed by merging one or more true clusters from the original data. Consequently, the number of possible estimated membership label configurations in underfitting cases is significantly reduced. This simplification provides considerable convenience for our subsequent proofs, as it narrows the range of scenarios that need to be considered.

As introduced in Jin et al. (2022), the NSP is a challenging property to establish. However, they provide a "stronger version of the k-means theorem," specifically Theorem 4.1 in their paper, which is highly useful. Leveraging this theorem, it is not difficult to show that the spectral clustering method described in Section 2 satisfies the NSP with high probability. In fact, this result is not limited to the specific case considered here; it can be similarly proven for many low-rank models with exact membership. The key lies in the property that the left singular vector matrices of the data matrix and the signal matrix differ by at most an orthogonal matrix, with their discrepancy being small. This property has been extensively studied, and for cases where the data matrix is bounded, the needed conclusion can almost directly follow by combining Theorem 4.1 in Jin et al. (2022) and Theorem 4.4 in Chen et al. (2021). Specifically, in our case, we have the following theorem, where  $\widehat{\mathbf{Z}}^{(m)}$  represents the membership labels obtained from spectral clustering:

**Theorem 1.** With probability at least  $1 - O(n^{-5})$ , for  $\forall 1 < m \leq K$ ,  $\mathbf{Z} \preceq \widehat{\mathbf{Z}}^{(m)}$  up to a permutation in the columns.

This result enables us to establish a theoretical guarantee for the consistency of our method under the following mild assumptions.

**Assumption 4.** Each entry in the signal matrix  $\mathbf{P}$  is bounded by a constant  $C_P$ , i.e.,  $|\mathbf{P}_{i,j}| \leq C_P$  for all  $i$  and  $j$ .

**Assumption 5.** As  $n$  and  $p$  increase, the ratio  $\frac{n}{N}$  is uniformly bounded both below and above by constants  $C_1$  and  $C_2$ , i.e.,  $C_1 \leq \frac{n}{N} \leq C_2$  for some constants  $C_1 > 0$  and  $C_2 > 0$ .

The first assumption is mild as it simply assumes that each entry of the cluster centers is bounded, which is a reasonable condition in most practical settings. The second assumption is also quite flexible, as it only requires that  $n$  and  $p$  grow at comparable rates, with their ratio being bounded both below and above. This condition ensures that  $n$  and  $p$  neither grow disproportionately large nor deviate significantly in scale.

Under the aforementioned assumptions, we now present the following theoretical guarantees for our method. In the theorems below, the statistic  $\phi_N^{(m)}$  is constructed as described in Algorithm 1. As previously mentioned, if we establish the result for this construction, then  $\phi_N^{(m)}$  defined in Algorithm 2 will also satisfy the same asymptotic properties.

**Theorem 2** (Null case:  $m = K$ ). Fix  $0 < \alpha < 1$ . As  $n \rightarrow \infty$ , we have:

$$\mathbb{P}(\phi_n^{(K)} \leq z_\alpha) \geq 1 - \alpha + o(1),$$

and

$$\mathbb{P}(\widehat{K}_\alpha \leq K) \geq 1 - \alpha + o(1).$$

It follows that  $\widehat{K}_\alpha^*$  is a level- $(1-\alpha)$  confidence lower level for  $K$ .

**Theorem 3** (Underfitting case:  $m < K$ ). *Fix  $0 < \alpha < 1$ . As  $N \rightarrow \infty$ , we have:*

$$\min_{1 \leq m < K} \{\phi_N^{(m)}\} \rightarrow \infty \quad \text{in probability},$$

and

$$\mathbb{P}(\widehat{K}_\alpha \neq K) \leq \alpha + o(1).$$

Now if we let  $\alpha$  depend on  $N$  and tend to 0 slowly enough, then we have proved  $\mathbb{P}(\widehat{K}_\alpha^* = K) \rightarrow 1$ . These results establish the asymptotic consistency of our method for both the null and underfitting cases, providing theoretical guarantees for its performance as  $N \rightarrow \infty$ .

**Remark 3.** *From our subsequent theorem proofs, it becomes evident that while we follow the formulation in Jin et al. (2022) by writing  $z_\alpha$ , this term can, in practice, be replaced by any constant. In principle,  $z_\alpha$  can also serve as a tuning parameter, offering additional flexibility in the application of our method.*

## 4 Numerical Studies

### 4.1 Theoretical Verification

We begin with numerical experiments to validate the asymptotic properties established in Theorems 2 and 3, particularly focusing on the theoretical constraint  $\beta \in (4, 8)$ . Our simulations aim to demonstrate that only within this range does the test statistic exhibit the desired asymptotic behavior: diverging to infinity in underfitting cases while converging to zero under the null hypothesis as  $N \rightarrow \infty$ .

To comprehensively evaluate the theoretical findings, we conduct extensive simulation studies across diverse model frameworks. The experiments encompass three distinct noise distributions—Bernoulli, Gamma, Gaussian, and Poisson—enabling us to assess the method’s robustness and consistency under varied stochastic conditions.

We implement a mixture model with  $K = 4$  components in high-dimensional settings. The data generation process consists of two phases: signal generation and noise incorporation.

In the signal generation phase, we first generate  $K$  cluster centers uniformly from the hypercube  $[-10, 10]^p$ . Each sample is then assigned to one of the  $K$  clusters, following a discrete uniform distribution with a probability of  $1/K$  for each cluster.

In the noise incorporation phase, the noise component is introduced according to one of four scenarios:

- **Bernoulli Noise Setting:** For each entry  $(i, j)$ , we add centered Bernoulli noise  $\epsilon_{ij} = B(p_{ij}) - p_{ij}$ , where  $p_{ij}$  is uniformly sampled from  $[0.1, 0.9]$ . This introduces binary, bounded perturbations with a maximum theoretical variance of 0.25.
- **Gamma Noise Setting:** For each entry  $(i, j)$ , we add centered Gamma noise  $\epsilon_{ij} = \Gamma(a_{ij}, b_{ij}) - a_{ij}b_{ij}$ , where  $a_{ij}$  and  $b_{ij}$  are chosen such that their product is uniformly sampled from  $[0.5, 5]$ . This introduces asymmetric perturbations with a maximum theoretical variance of 125.
- **Gaussian Noise Setting:** For each entry  $(i, j)$ , we introduce centered Gaussian noise  $\epsilon_{ij} \sim N(0, \sigma_{ij})$ , where  $\sigma_{ij}$  is uniformly sampled from  $[1, 100]$ . This represents continuous, symmetric perturbations with heteroscedastic variance.

- **Poisson Noise Setting:** We incorporate centered Poisson noise  $\epsilon_{ij} = P(\lambda_{ij}) - \lambda_{ij}$  for each entry  $(i, j)$ , where the intensity parameter  $\lambda_{ij}$  is uniformly sampled from  $\{1, \dots, 100\}$ . This generates discrete, asymmetric perturbations with heavy-tailed characteristics.

The denominator of the test score takes the form  $\sqrt{\sigma^8 N^\beta}$ , where both  $\sigma$  and  $\beta$  play crucial roles in its asymptotic behavior. For  $\sigma$ , we employ the theoretical upper bounds of the respective noise variances:  $\sigma = 100$  for both Poisson and Gaussian settings (corresponding to their maximum variances),  $\sigma = 0.25$  for the Bernoulli setting (matching its theoretical maximum variance), and  $\sigma = 125$  for the Gamma setting (matching its theoretical maximum variance).

To investigate the impact of  $\beta$ , we examine five values:  $\beta = 4, 5, 6, 7, 8$ . For each combination of  $\beta$  and noise type, we analyze the behavior of the test statistic across varying sample sizes, with dimension  $p = 10n$  for  $n \in \{50, 100, \dots, 1000\}$ . To ensure stability, each configuration is replicated 100 times, and we present the averaged test scores.

The simulation results are presented in Figures 1–4. For each noise type and each value of  $\beta$ , we use dual y-axes to display all test statistics in one plot: the left y-axis (black) corresponds to the underfitting cases ( $m < K$ ), while the right y-axis (purple) corresponds to the null case ( $m = K$ ).

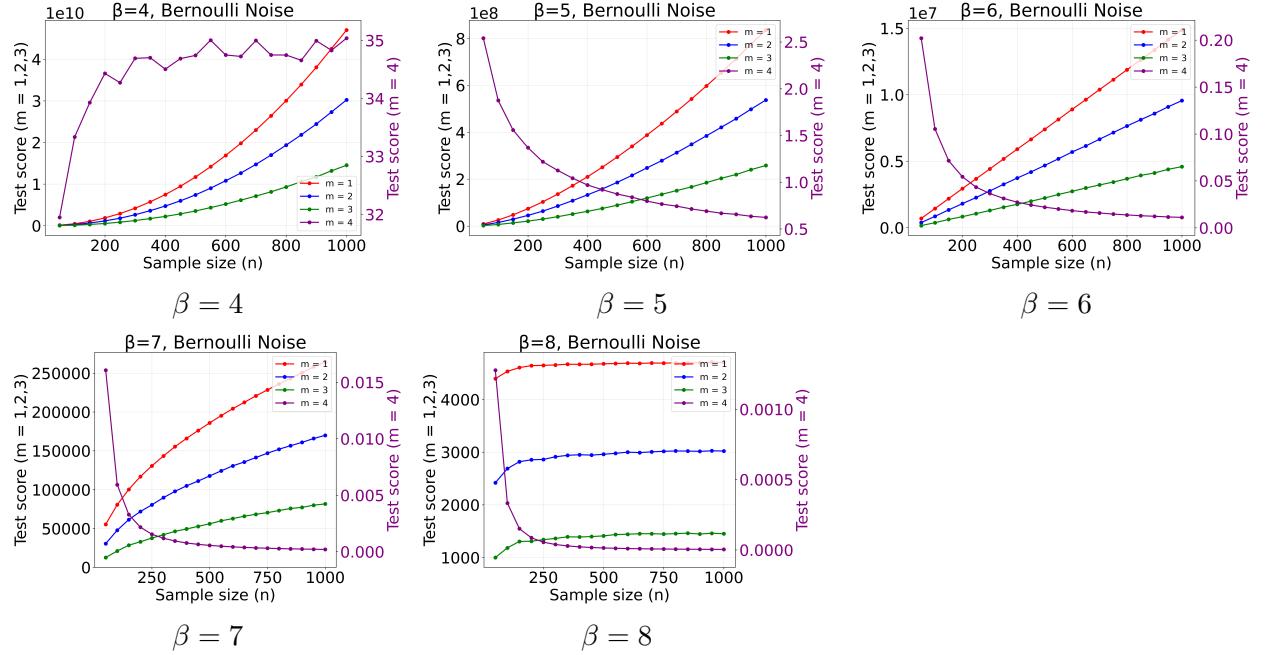


Figure 1: Test scores under Bernoulli noise for different values of  $\beta$ . For each subplot, the left y-axis (black) corresponds to underfitting cases ( $m = 1, 2, 3$ ), while the right y-axis (purple) corresponds to the correct specification ( $m = 4$ ).

The results demonstrate remarkable consistency across all three noise types. For  $\beta = 5, 6, 7$ , in underfitting cases ( $m < K$ ), the test statistics diverge to infinity as the sample size increases; in contrast, for the null case ( $m = K$ ), the test score converges to 0 as the sample size increases. This pattern holds regardless of the noise distribution, demonstrating the robustness of our method. However, this pattern breaks down at the boundary cases: when  $\beta = 4$ , the test scores fail to converge to 0 in the null case, and when  $\beta = 8$ , the test scores fail to diverge to infinity in the underfitting cases, as expected.

**Remark 4.** In fact, in certain scenarios (e.g., when the noise within the same cluster is identically

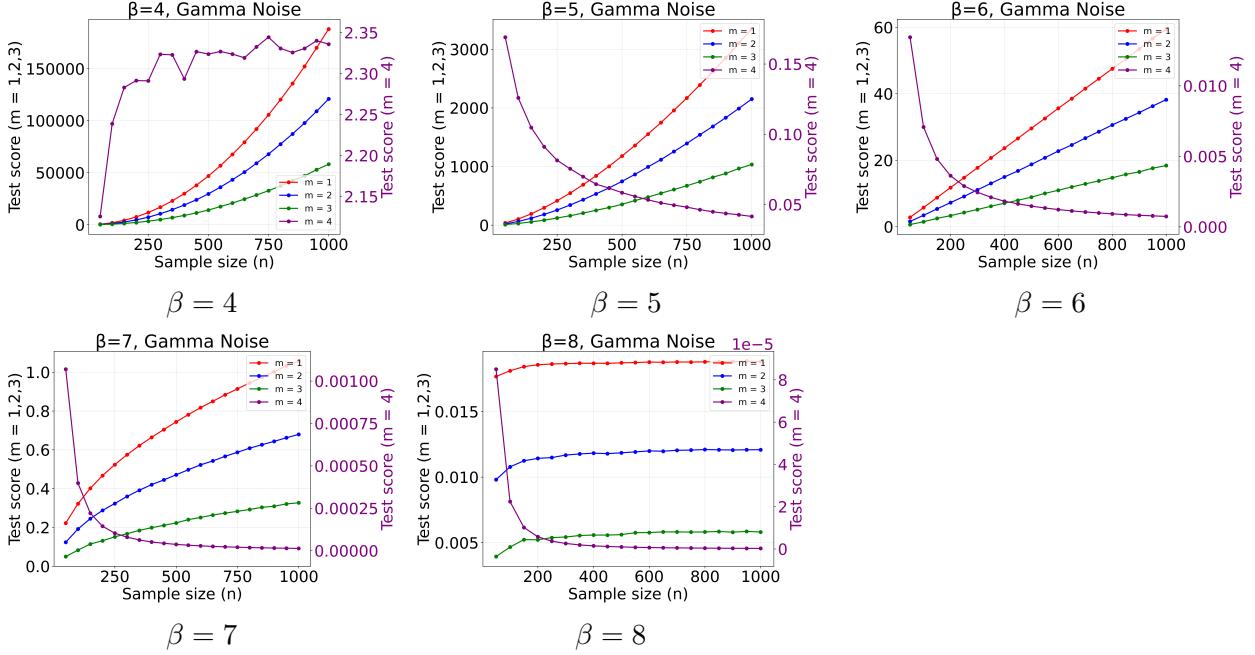


Figure 2: Test scores under Gamma noise for different values of  $\beta$ . For each subplot, the left y-axis (black) corresponds to underfitting cases ( $m = 1, 2, 3$ ), while the right y-axis (purple) corresponds to the correct specification ( $m = 4$ ).

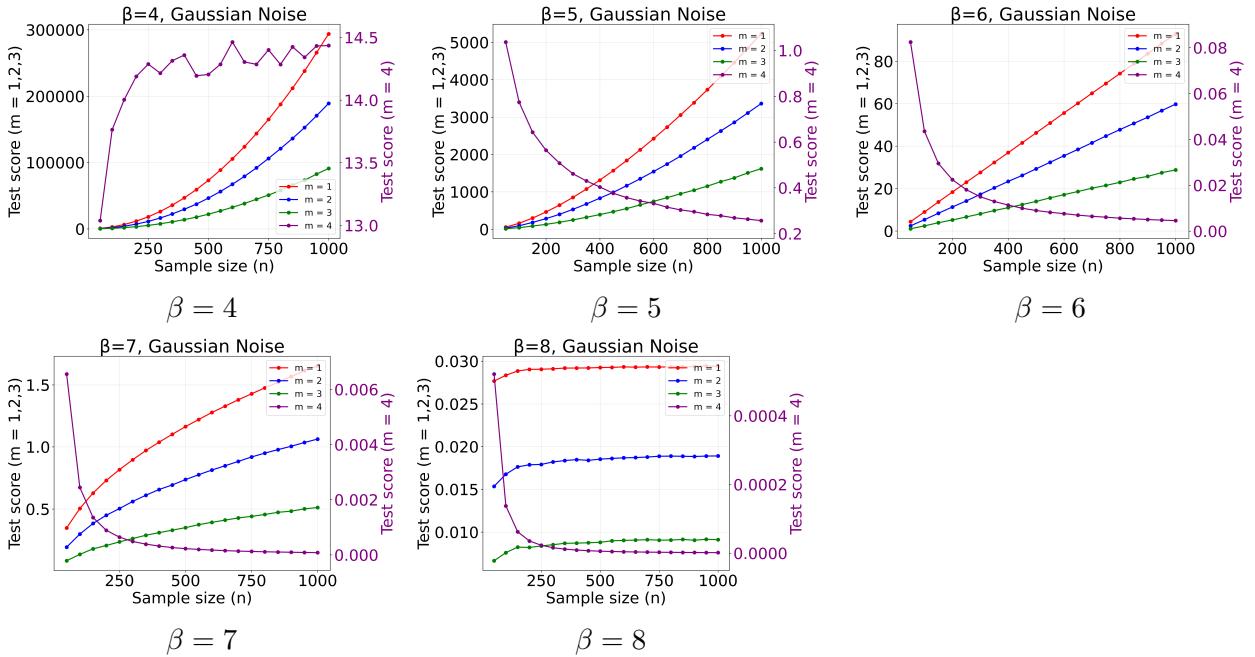


Figure 3: Test scores under Gaussian noise for different values of  $\beta$ . For each subplot, the left y-axis (black) corresponds to underfitting cases ( $m = 1, 2, 3$ ), while the right y-axis (purple) corresponds to the correct specification ( $m = 4$ ).

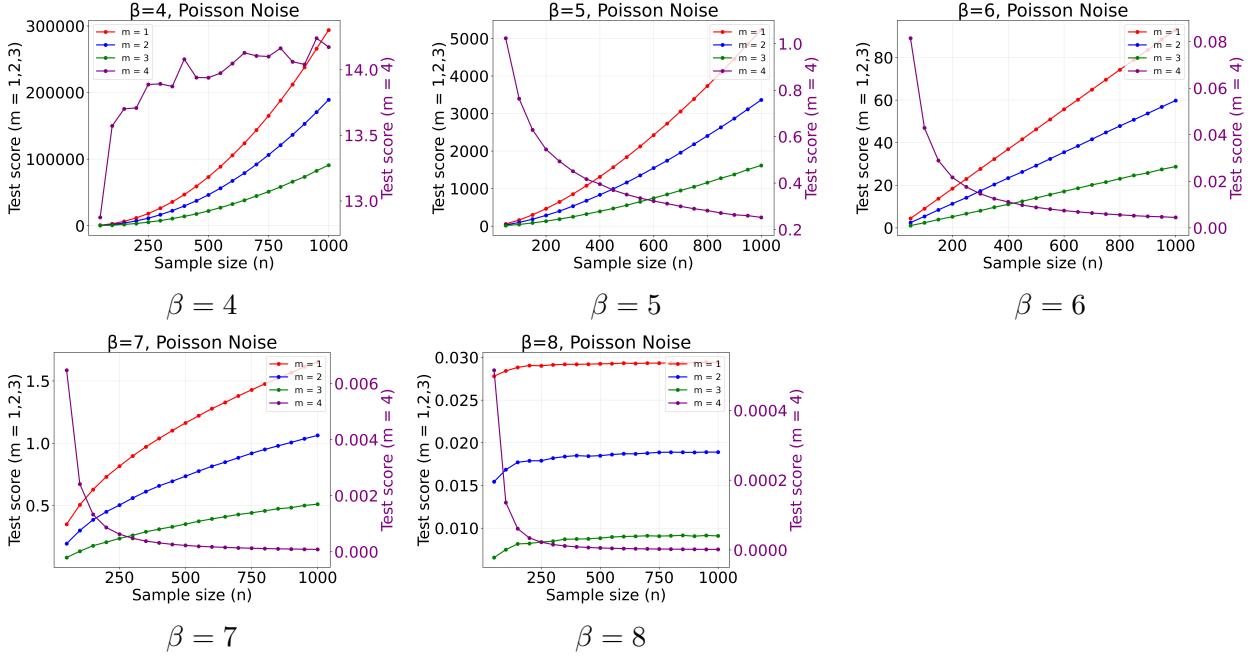


Figure 4: Test scores under Poisson noise for different values of  $\beta$ . For each subplot, the left y-axis (black) corresponds to underfitting cases ( $m = 1, 2, 3$ ), while the right y-axis (purple) corresponds to the correct specification ( $m = 4$ ).

*distributed), we have also observed that under the null case, the test score decreases as  $n$  increases when  $\beta = 4$ . However, given the counterexamples provided above, it may be unlikely to establish, in general, conclusions such as Theorems 2 and 3 when  $\beta = 4$  or  $\beta = 8$ .*

## 4.2 Performance Comparison

We evaluate the accuracy of our proposed method compared to the Bayesian Information Criterion (BIC) method under varying values of  $n$ ,  $\Delta$ , and  $n/p$ . We also attempted to incorporate the method described in Manole and Khalili (2021) into the comparison by utilizing their publicly available R package, **GroupSortFuse**. The package supports Poisson mixture models, but it is limited to one-dimensional settings and cannot handle higher-dimensional data. For Gaussian mixture models, the package can only handle models where the variances are the same, which is a significant limitation for our experiments. Due to these constraints, we ultimately chose not to include this method in our benchmarks.<sup>1</sup>

We ultimately selected BIC as the benchmark for the comparisons presented in this section. To compute BIC, we used the **flexmix** package in R, which supports only high-dimensional Gaussian and Poisson mixture models. As a result, we limited our comparisons to these two types of noise.

In this subsection, the StGoF method we employed follows Algorithm 2, where we estimate  $\sigma$  directly, assuming no prior knowledge of  $\sigma$ . Additionally, for both BIC and StGoF, we used the same upper bound for  $K$  to ensure a fair comparison of computational efficiency in Section 4.2.4.

<sup>1</sup>Due to time constraints, we were unable to fine-tune the method, but future comparisons may address this issue.

### 4.2.1 Gaussian Mixture Models

In this experiment, we fixed  $K = 3$  and considered nine combinations of  $(n, p)$  values. These combinations were grouped into three categories based on the ratio  $\frac{n}{p} = \frac{1}{5}, 1, 5$ . For each  $(n, p)$  pair, we further examined five different values of the minimum distance between cluster centers  $\Delta$ . For each combination of  $(n, p, \Delta)$ , we first generated  $K$  cluster centers uniformly within the hypercube  $[200, 400]^p$ . To ensure that the minimum distance between cluster centers was exactly  $\Delta$ , the initial centers were scaled proportionally. Each sample was then assigned to one of the  $K$  clusters, following a discrete uniform distribution where each cluster was selected with a probability of  $1/K$ .

After generating the cluster centers, we added independent noise drawn from  $N(0, 1)$  to each component of the cluster centers, resulting in a data matrix. This process was repeated 100 times. For each of the 100 datasets, we applied both the BIC and StGoF methods to estimate  $K$  and computed the accuracy of the estimates. The final results are presented in Figure 5.

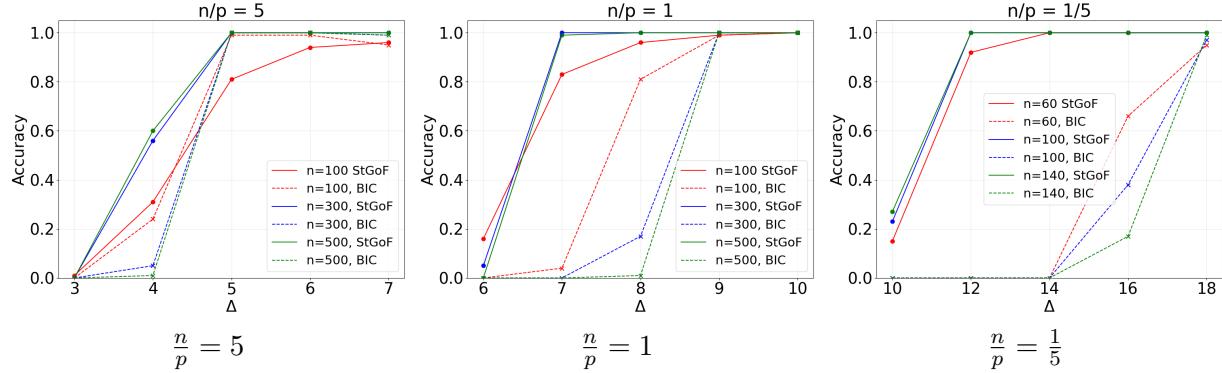


Figure 5: Performance comparison of StGoF and BIC in Gaussian mixture models: Solid lines represent StGoF, and dashed lines represent BIC. The accuracy is plotted against  $\Delta$  for different values of  $n$ .

From the figure, it can be observed that StGoF consistently outperforms BIC in most settings. This advantage is particularly pronounced in scenarios with smaller  $\frac{n}{p}$  ratios, where StGoF successfully predicts  $K$  even under more challenging conditions where BIC fails to provide accurate estimates.

### 4.2.2 Poisson Mixture Models

Similar to the Gaussian mixture models described in the previous section, we fixed  $K = 7$  and considered nine combinations of  $(n, p)$  values, grouped into three categories based on the ratio  $\frac{n}{p} = \frac{1}{5}, 1, 5$ . For each  $(n, p)$  pair, we further examined five different values of the minimum distance between cluster centers  $\Delta$ .

In the case of Poisson mixture models, each cluster center corresponds to a  $p$ -dimensional rate parameter vector for a Poisson distribution. Specifically, for each of the  $K$  clusters, we first generated a  $p$ -dimensional vector of rate parameters uniformly within the hypercube  $[1, 10]^p$ . To ensure that the minimum distance between the cluster centers was exactly  $\Delta$ , the initially generated cluster centers were then scaled proportionally. After scaling the cluster centers, each sample was assigned to one of the  $K$  clusters with a probability of  $1/K$ , and data points for each cluster were sampled independently from a Poisson distribution in each dimension, with the corresponding rate parameter for that dimension.

After generating the data, we applied both the BIC and StGoF methods to estimate  $K$  for each of the 100 generated datasets and computed the accuracy of the estimates. The final results are presented in Figure 6, where we compare the performance of both methods under different experimental settings.

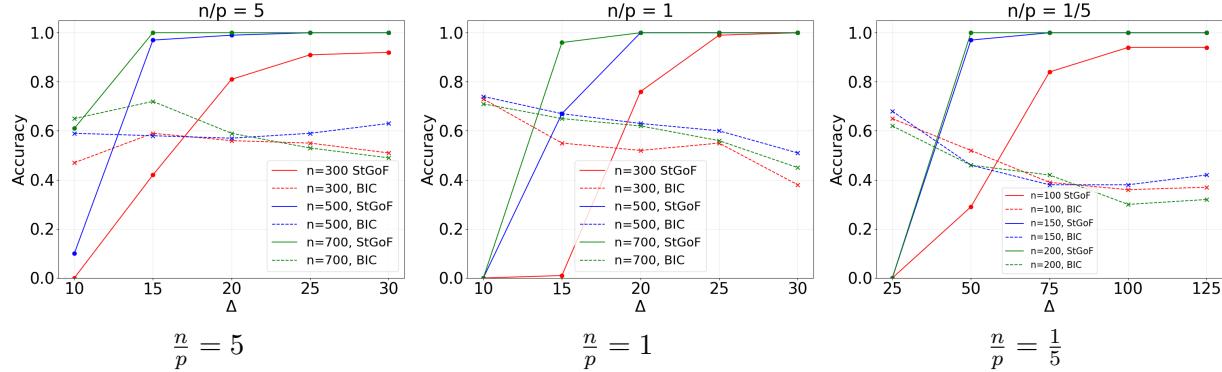


Figure 6: Performance comparison of StGoF and BIC in Poisson mixture models: Solid lines represent StGoF, and dashed lines represent BIC. The accuracy is plotted against  $\Delta$  for different values of  $n$ .

As can be seen, the accuracy of StGoF improves rapidly with an increase in  $\Delta$ . In contrast, the performance of BIC using R’s `flexmix` package is less satisfactory, especially for relatively larger values of  $K$ .<sup>2</sup>

#### 4.2.3 Multi-distribution Mixture Models

In this experiment, we follow a similar data generation process as described in the previous sections, but with a few key differences. The combinations of  $(n, p)$  values are divided into three groups based on the ratio  $n/p = 1/5, 1, 5$ , resulting in a total of nine distinct combinations. For each combination of  $(n, p)$ , five different values of the minimum distance between cluster centers ( $\Delta$ ) are considered. Each of the  $(n, p, \Delta)$  combinations is then repeated 100 times to compute the accuracy of the clustering method.

In this case, we set  $K = 7$  clusters. The process of generating the cluster centers begins by first selecting  $K$  cluster centers within the hypercube  $[-10, 10]^p$ . These initial centers are then scaled to ensure that the minimum distance between them is exactly  $\Delta$ .

The data generation for each sample proceeds as follows. First, for each cluster, a corresponding  $p$ -dimensional vector of noise parameters is generated from a uniform distribution between 1 and 10. This vector of noise parameters corresponds to the variance for each cluster in each dimension.

For each sample, the  $p$  dimensions are randomly split into two groups:  $p/2$  dimensions will be assigned Poisson noise, and the remaining  $p/2$  dimensions will be assigned Gaussian noise. For the  $p/2$  dimensions assigned to Poisson noise, the data points are sampled from a Poisson distribution with the corresponding rate parameters drawn from the noise parameter vector for that cluster. For the remaining  $p/2$  dimensions, Gaussian noise is applied, where the data points are sampled from a normal distribution with zero mean and variance determined by the corresponding noise parameter vector.

In our experiments, we compare the clustering results obtained from our proposed method with those from the gap statistics approach (Tibshirani et al., 2000). We applied the `ClusGap` function

<sup>2</sup>Based on my tests, BIC performs reasonably well for smaller values of  $K$ .

from the R package, using `kmeans` as the clustering algorithm and setting the number of bootstrap iterations to 100 ( $B=100$ ). The gap statistic was calculated for different numbers of clusters, and the optimal number of clusters was determined based on the gap statistic values.

**Remark 5.** Utilizing PAM instead of  $k$ -means can improve clustering accuracy; however, this comes at the expense of computational efficiency.

The results from this method are presented in Figure 7. In most cases, our method outperforms the gap statistics approach, especially as the sample size  $n$  and parameter  $\Delta$  increase.

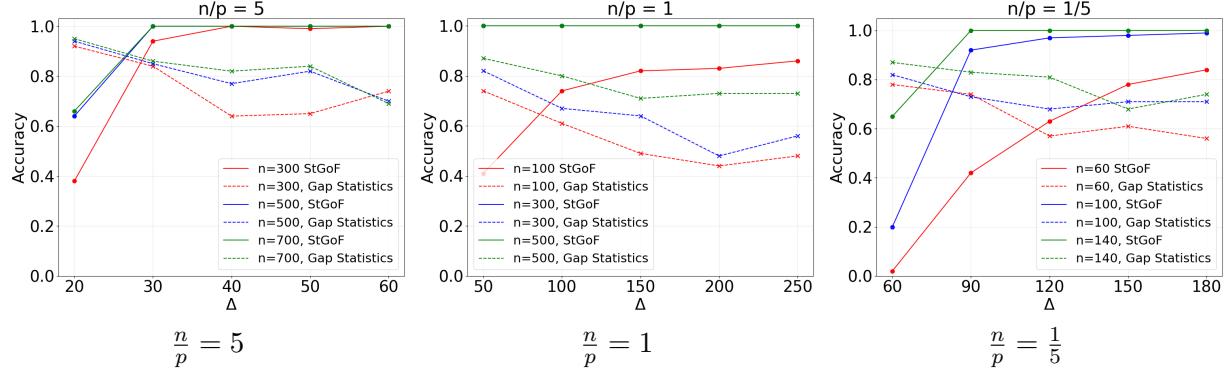


Figure 7: Performance comparison of StGoF and gap statistics approach in Multi-distribution mixture models: Solid lines represent StGoF, and dashed lines represent gap statistics approach. The accuracy is plotted against  $\Delta$  for different values of  $n$ .

#### 4.2.4 Computational Efficiency

In this subsection, we evaluate the computational efficiency of our proposed method compared to BIC, based on the experiments described in Section 4.2.1. The computational time for each combination of  $(n, p, \Delta)$  was averaged over 100 runs. To streamline the presentation, we selected three representative values of  $\Delta$  for each pair  $(n, p)$ . The results are visualized in Figure 8 as logarithmic computational time plotted against  $n$ . From the figure, it is evident that our method is significantly faster than BIC, often achieving runtimes that are approximately two orders of magnitude smaller. This stark contrast highlights the efficiency and practicality of our approach in these scenarios.

### 4.3 Real Data Analysis

In this subsection, we apply our method to the United States 112th Senate Roll Call Votes dataset, which records the voting behavior of U.S. senators over  $J = 486$  roll calls. Following the preprocessing steps described in Lyu et al. (2024), we encode the original categorical voting data into binary responses and remove senators with excessive missing votes or those not affiliated with the two major political parties. For the remaining  $N = 94$  senators, missing entries are imputed probabilistically based on their individual voting patterns.

From our theoretical analysis, it is evident that the guarantees of our method rely on the assumption of exact recovery, which requires a sufficiently large sample size. In small-sample scenarios, our method, being significantly faster than traditional likelihood-based approaches, can serve as a rough estimator for the number of components and provide a practical initialization

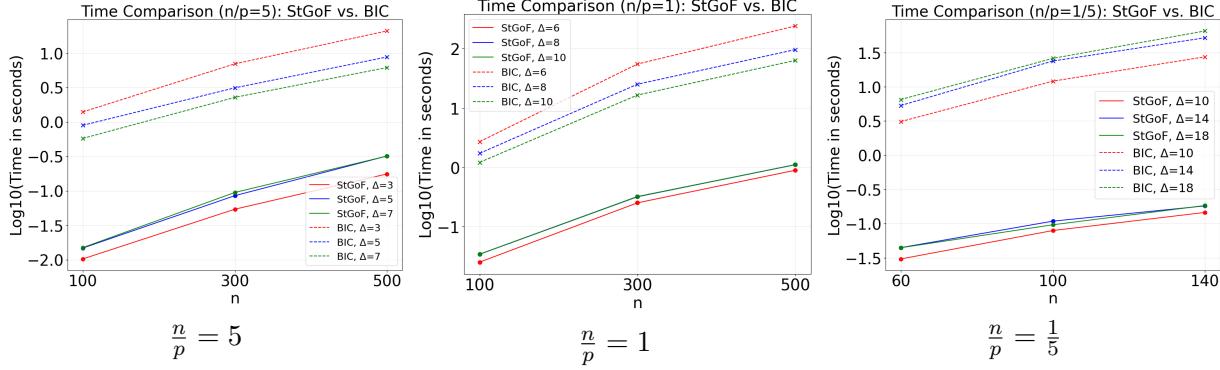


Figure 8: Logarithmic computational time as a function of  $n$  for three values of  $\Delta$ . The comparison demonstrates the superior efficiency of our proposed method (solid lines) over BIC (dashed lines). These results correspond to the experiments detailed in Section 4.2.1.

for more refined methods. However, the challenge of estimating  $\sigma$  in such settings necessitates adjustments to enhance robustness.

To address this, we modify the definition of  $C_N$  to  $C_N = 2n^{\beta/2}p^{\beta/2}$ , with  $\beta$  increased to 6, to compensate for the lack of a reliable  $\sigma$  estimate. Applying this modification to the Senate Roll Call Votes data, our method estimates the number of components to be  $\hat{K} = 2$ . This result aligns with the known political structure of the dataset, which comprises two major political parties: Democrats and Republicans. This demonstrates the practical utility of our method in real-world scenarios, particularly as a fast and interpretable tool for exploratory data analysis.

## 5 Discussion

The cut-off value of  $z_\alpha$  in our framework is theoretically grounded, but in practice, it can be replaced by any constant, depending on the specific context or problem setting. While increasing  $n$  in simulation studies helps mitigate issues related to this cut-off, determining the appropriate value for real-world data remains challenging. This highlights the need for more adaptive frameworks that can better account for varying data structures and problem-specific requirements.

A major difficulty in real-world applications arises from the assumption that data is generated by adding noise to predefined cluster centers with zero mean. This assumption works in controlled settings but is overly restrictive for many real-world datasets where cluster centers are not well-defined. As a result, the test scores may not behave as expected under the null hypothesis, necessitating more flexible approaches that account for inherent data variability.

The estimation of the variance parameter  $\sigma$  is another crucial factor influencing the reliability of test scores. Although  $\sigma$  is treated as constant in our model, real-world data often deviates from this assumption. Inaccurate estimation of  $\sigma$  can significantly distort the results, and its estimation typically requires prior knowledge of the clustering structure, which is often unavailable. Future work could focus on developing methods for robustly estimating  $\sigma$  without relying on prior clustering knowledge. Such methods would enhance the flexibility of our approach and improve its performance in real-world scenarios where the data structure is less predictable. Nevertheless, given the significant computational advantage of this method, even in its current form, it can still serve as a reliable initialization for the number of components, even without further modifications.

## References

- Abbe, E., Fan, J., and Wang, K. (2022). An  $\ell_p$  theory of pca and spectral clustering. *The Annals of Statistics*, 50(4):2359 – 2385.
- Ahmad, A. and Dey, L. (2011). A k-means type clustering algorithm for subspace clustering of mixed numeric and categorical datasets. *Pattern Recognit. Lett.*, 32(7):1062–1069.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723.
- Bechtel, Y. C., Bonaiti-Pellie, C., Poisson, N., Magnette, J., and Bechtel, P. (1993). A population and family study n-acetyltransferase using caffeine urinary metabolites. *Clinical Pharmacology & Therapeutics*, 54.
- Biernacki, C., Celeux, G., and Govaert, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22:719 – 725.
- Bouguila, N. and Fan, W. (2020). *Mixture Models and Applications*. Springer.
- Cai, T. T. and Zhang, A. (2018). Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *The Annals of Statistics*, 46:60–89.
- Chen, J. and Kalbfleisch, J. D. (1996). Penalized minimum-distance estimates in finite mixture models. *Canadian Journal of Statistics-revue Canadienne De Statistique*, 24:167–175.
- Chen, J. and Li, P. (2009). Hypothesis test for normal mixture models: The em approach. *The Annals of Statistics*, 37(5A):2523 – 2542.
- Chen, X. and Yang, Y. (2021). Cutoff for exact recovery of gaussian mixture models. *IEEE Transactions on Information Theory*, 67:4223–4238.
- Chen, Y., Chi, Y., Fan, J., and Ma, C. (2021). Spectral methods for data science: A statistical perspective. *Foundations and Trends in Machine Learning*, 14(5):566–806.
- Dacunha-Castelle, D. and Gassiat, E. (1999). Testing the order of a model using locally conic parametrization: Population mixtures and stationary arma processes. *The Annals of Statistics*, 27(4):1178–1209.
- Dreveton, M., Gözeten, A., Grossglauser, M., and Thiran, P. (2024). Universal lower bounds and optimal rates: Achieving minimax clustering error in sub-exponential mixture models. In *Proceedings of the Thirty Seventh Conference on Learning Theory*, volume 247 of *Proceedings of Machine Learning Research*, pages 1451–1485. PMLR.
- Drton, M. and Plummer, M. (2013). A bayesian information criterion for singular models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79.
- Dudoit, S. and Fridlyand, J. (2002). A prediction-based resampling method for estimating the number of clusters in a dataset. *Genome Biology*, 3.
- Foss, A. H., Markatou, M., and Ray, B. (2019). Distance metrics and clustering methods for mixed-type data. *International Statistical Review*, 87(1):80–109.

- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT press Cambridge.
- Gower, J. C. (1971). A general coefficient of similarity and some of its properties. *Biometrics*, 27(4):857–871.
- Heinrich, P. and Kahn, J. (2018). Strong identifiability and optimal minimax rates for finite mixture estimation. *The Annals of Statistics*, 46(6A):2844 – 2870.
- Hennig, C. and Liao, T. F. (2013). How to find an appropriate clustering for mixed-type variables with application to socio-economic stratification. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 62(3):309–369.
- Hennig, C. and Lin, C.-J. (2015). Flexible parametric bootstrap for testing homogeneity against clustering and assessing the number of clusters. *Statistics and Computing*, 25(4):821–833.
- Horn, R. A. and Johnson, C. R. (1985). *Matrix Analysis*. Cambridge University Press.
- Huang, Z. (1998). Extensions to the k-means algorithm for clustering large data sets with categorical values. *Data Mining and Knowledge Discovery*, 2(3):283–304.
- James, L. F., Marchette, D. J., and Priebe, C. E. (2001). Consistent estimation of mixture complexity. *The Annals of Statistics*, 29(5):1281 – 1296.
- Jin, J., Ke, Z., and Luo, S. (2018). Network global testing by counting graphlets. *Proceedings of the 35th International Conference on Machine Learning*, 80:2333–2341.
- Jin, J., Ke, Z. T., Luo, S., and Wang, M. (2022). Optimal estimation of the number of communities. *Journal of the American Statistical Association*, 118(543):2101–2116.
- Keribin, C. (2000). Consistent estimation of the order of mixture models. *The Indian Journal of Statistics*, 62:49–66.
- Leroux, B. G. (1992). Consistent estimation of a mixing distribution. *The Annals of Statistics*, 20(3):1350 – 1360.
- Li, P. and Chen, J. (2010). Testing the order of a finite mixture. *Journal of the American Statistical Association*, 105(491):1084–1092.
- Liu, X. and Shao, Y. (2003). Asymptotics for likelihood ratio tests under loss of identifiability. *The Annals of Statistics*, 31(3):807 – 832.
- Lyu, Z., Chen, L., and Gu, Y. (2024). Degree-heterogeneous latent class analysis for high-dimensional discrete data. *arXiv preprint arXiv:2402.18745*.
- Manole, T. and Khalili, A. (2021). Estimating the number of components in finite mixture models via the group-sort-fuse procedure. *The Annals of Statistics*, 49(6):3043 – 3069.
- McLachlan, G. J. (1987). On bootstrapping the likelihood ratio test statistic for the number of components in a normal mixture. *Applied statistics*, 36:318–324.
- McLachlan, G. J., Lee, S. X., and Rathnayake, S. I. (2019). Finite mixture models. *Annual Review of Statistics and Its Application*, 6(Volume 6, 2019):355–378.
- McLachlan, G. J. and Peel, D. (2004). *Finite Mixture Models*. John Wiley & Sons.

- Mixon, D. G., Villar, S., and Ward, R. (2017). Clustering subgaussian mixtures by semidefinite programming. *Information and Inference: A Journal of the IMA*, 6:389–415.
- Ndaoud, M. (2018). Sharp optimal recovery in the two-component gaussian mixture model. *arXiv preprint arXiv:1812.08078*.
- Schlattmann, P. (2009). *Medical Applications of Finite Mixture Models*. Springer Science & Business Media.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6:461–464.
- Srivastava, P., Sarkar, P., and Hanusanto, G. A. (2019). A robust spectral clustering algorithm for sub-gaussian mixture models with outliers. *arXiv preprint arXiv:1912.07546*.
- Tibshirani, R. and Walther, G. (2005). Cluster validation by prediction strength. *Journal of Computational and Graphical Statistics*, 14(3):511–528.
- Tibshirani, R., Walther, G., and Hastie, T. J. (2000). Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63.
- Vempala, S. and Wang, G. (2004). A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68:841–860.
- Volkovich, Z., Barzily, Z., Weber, G.-W., Toledano-Kitai, D., and Avros, R. (2011). Resampling approach for cluster model selection. *Machine Learning*, 85:209–248.
- Wang, C. and Yang, Y. (2024). Estimating the number of components in finite mixture models via variational approximation. *arXiv preprint arXiv:2404.16746*.
- Woo, M.-J. and Sriram, T. N. (2006). Robust estimation of mixture complexity. *Journal of the American Statistical Association*, 101:1475 – 1486.
- Zhang, A. Y. and Zhou, H. H. (2021). Optimality of spectral clustering for gaussian mixture model. *The Annals of Statistics*, 49(5):2506–2530.
- Zhang, A. Y. and Zhou, H. H. (2022). Leave-one-out singular subspace perturbation analysis for spectral clustering. *arXiv preprint arXiv:2205.14855*.

## Supplementary Material

### S.1 Proof of Theorem 2

The second claim follows directly from the first claim. We will focus on the first claim.

As a natural corollary of Theorem 1, we already know

$$\mathbb{P}(\widehat{\mathbf{Z}}^{(K)} \neq \mathbf{Z}) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Comparing to  $\widehat{\mathbf{P}}^{(K)} = \widehat{\mathbf{Z}}^{(K)}(\mathbf{Z}^\top \mathbf{Z})^{-1}(\widehat{\mathbf{Z}}^{(K)})^\top \mathbf{X}$ , we introduce the proxies of  $\widehat{\mathbf{P}}^{(K)}$ ,  $Q_N^{(K,0)}$  and  $\phi_N^{(K,0)}$  below:

$$\widehat{\mathbf{P}}^{(K,0)} = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1}\mathbf{Z}^\top \mathbf{X},$$

$$Q_N^{(K,0)} = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} (\mathcal{S}(\mathbf{A})_{i_1, i_2} - (\mathcal{S}(\widehat{\mathbf{P}}^{(K,0)})_{i_1, i_2})(\mathcal{S}(\mathbf{A})_{i_2, i_3} - (\mathcal{S}(\widehat{\mathbf{P}}^{(K,0)})_{i_2, i_3})(\mathcal{S}(\mathbf{A})_{i_3, i_4} - (\mathcal{S}(\widehat{\mathbf{P}}^{(K,0)})_{i_3, i_4})) \\ (\mathcal{S}(\mathbf{A})_{i_4, i_1} - (\mathcal{S}(\widehat{\mathbf{P}}^{(K,0)})_{i_4, i_1}),$$

$$\phi_N^{(K,0)} = Q_N^{(K,0)} / \sqrt{C_N}.$$

For fixed  $t \in \mathbb{R}$ ,  $|\mathbb{P}(\phi_N^{(K)} \leq t) - \mathbb{P}(\phi_N^{(K,0)} \leq t)| \leq \mathbb{P}(\widehat{\mathbf{Z}}^{(K)} \neq \mathbf{Z}) \rightarrow 0$  as  $N \rightarrow \infty$ .

Next we will show: for fixed  $t \in \mathbb{R}$ ,  $\mathbb{P}(\phi_N^{(K,0)} \geq z_\alpha) \geq 1 - \alpha + o(1)$  as  $n \rightarrow \infty$ . Hence  $\mathbb{P}(\phi_N^{(K)} \geq z_\alpha) \geq 1 - \alpha + o(1)$  as  $n \rightarrow \infty$ .

We define  $\tilde{Q}_N = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$ , where  $\mathbf{W} = \mathcal{S}(\mathbf{A}) - \mathcal{S}(\mathbf{P})$ .

We have the following lemma, which will be proved later.

**Lemma S.1.**  $\tilde{Q}_N / \sqrt{\text{Var}(\tilde{Q}_N)} \rightarrow \mathbb{N}(0, 1)$  in law.

**Lemma S.2.**  $E[(Q_N^{(K,0)} - \tilde{Q}_N)^2] = O(N^4)$

If we admit lemmas above, then we can rewrite  $\phi_N^{(K,0)}$  as  $\frac{\tilde{Q}_N}{\sqrt{8C_N}} + \frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{8C_N}}$ . By Lemma S.2, we know  $\mathbb{E}[(\frac{Q_N^{(K,0)} - \tilde{Q}_N}{\sqrt{8C_N}})^2] = \frac{O(N^4)}{\sigma^8 N^\beta} \rightarrow 0$ , hence  $\frac{Q_N^{(K,0)} - \tilde{Q}_N}{\sqrt{8C_N}} \rightarrow 0$  in probability.

Therefore, for any fixed  $\epsilon$  such that  $0 < \epsilon < z_\alpha$ , we have

$$\begin{aligned} \mathbb{P}(\phi_N^{(K,0)} > z_\alpha) &= \mathbb{P}\left(\frac{\tilde{Q}_N}{\sqrt{C_N}} + \frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}} > z_\alpha\right) \\ &= \mathbb{P}\left(\frac{\tilde{Q}_N}{\sqrt{C_N}} + \frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}} > z_\alpha, \left|\frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}}\right| > \epsilon\right) + \mathbb{P}\left(\frac{\tilde{Q}_N}{\sqrt{C_N}} + \frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}} > z_\alpha, \left|\frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}}\right| \leq \epsilon\right) \\ &\leq \mathbb{P}\left(\left|\frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}}\right| > \epsilon\right) + \mathbb{P}\left(\frac{\tilde{Q}_N}{\sqrt{C_N}} > z_\alpha - \epsilon\right) \\ &\leq \mathbb{P}\left(\left|\frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}}\right| > \epsilon\right) + \mathbb{P}\left(\frac{\tilde{Q}_N}{\sqrt{\text{Var}(\tilde{Q}_N)}} > z_\alpha - \epsilon\right) \end{aligned}$$

where the last inequality follows from  $\text{Var}(\tilde{Q}_N) \leq C_N$ .

Then we can choose  $\epsilon = z_{s\alpha}$  where  $s$  is slightly larger than 1. Then we have

$$\mathbb{P}(\phi_N^{(K,0)} > z_\alpha) \leq \mathbb{P}\left(\left|\frac{(Q_N^{(K,0)} - \tilde{Q}_N)}{\sqrt{C_N}}\right| > s\alpha\right) + s\alpha = o(1) + s\alpha \quad \text{as } N \rightarrow \infty.$$

Let  $s$  approaches 1, we obtain that  $\mathbb{P}(\phi_N^{(K,0)} > z_\alpha) \leq \alpha + o(1)$ .

We now proceed to demonstrate that  $\text{Var}(\tilde{Q}_N) \leq C_N$ . Consider any ordered quadruple  $(i, j, k, \ell)$  with four distinct indices, there are 8 summands in the definition of  $\tilde{Q}_N$  whose values are exactly the same; these summands correspond to  $(i_1, i_2, i_3, i_4) \in \{(i, j, k, \ell), (j, k, \ell, i), (k, \ell, i, j), (\ell, i, j, k), (k, j, i, \ell), (j, i, \ell, k), (i, \ell, k, j), (\ell, k, j, i)\}$ , respectively. We treat these 8 summands as in an equivalent class. Denote by  $C(I_N)$  the collection of all such equivalent classes of four distinct nodes in  $\{1, \dots, N\}$ . Then, for any doubly indexed sequence  $\{x_{ij}\}_{1 \leq i \neq j \leq N}$  such that  $x_{ij} = x_{ji}$ , it is true that  $\sum_{i_1, i_2, i_3, i_4(\text{dist})} x_{i_1 i_2} x_{i_2 i_3} x_{i_3 i_4} x_{i_4 i_1} = 8 \sum_{C(I_N)} x_{i_1 i_2} x_{i_2 i_3} x_{i_3 i_4} x_{i_4 i_1}$ . In particular,

$$\tilde{Q}_N = 8 \sum_{C(I_n)} \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$$

. Since  $\mathbf{W}_{ij}$  and  $\mathbf{W}_{i'j'}$  are independent if  $(i, j) \neq (i', j')$  and  $E[\mathbf{W}_{ij}] = 0$ , the summands are uncorrelated of each other. As a result,

$$\begin{aligned} \text{Var}(\tilde{Q}_n) &= 64 \sum_{C(I_n)} \text{Var}(\mathbf{W}_{i_1 i_2}) \text{Var}(\mathbf{W}_{i_2 i_3}) \text{Var}(\mathbf{W}_{i_3 i_4}) \text{Var}(\mathbf{W}_{i_4 i_1}) \\ &= 8 \sum_{i_1, i_2, i_3, i_4(\text{dist})} \text{Var}(\mathbf{W}_{i_1 i_2}) \text{Var}(\mathbf{W}_{i_2 i_3}) \text{Var}(\mathbf{W}_{i_3 i_4}) \text{Var}(\mathbf{W}_{i_4 i_1}) \\ &= 8 \sum_{\substack{i_1, i_2, i_3, i_4(\text{dist}) \\ |\{i_1, i_2, i_3, i_4\} \cap \{1, \dots, n\}|=2}} \text{Var}(\mathbf{W}_{i_1 i_2}) \text{Var}(\mathbf{W}_{i_2 i_3}) \text{Var}(\mathbf{W}_{i_3 i_4}) \text{Var}(\mathbf{W}_{i_4 i_1}) \\ &\leq 2n^2 p^2 \sigma^8 \\ &= C_N. \end{aligned}$$

## S.2 Proof of Lemma S.1

For  $1 \leq M \leq N$ , define the  $\sigma$ -algebra  $\mathcal{F}_{N,M} = \sigma(\{\mathcal{S}(\mathbf{X})_{ij} : 1 \leq i < j \leq M\})$  and

$$Y_{N,M} = S_{N,M} - S_{N,M-1},$$

where  $S_{N,0} = 0$  and

$$S_{N,M} = \frac{\sum_{(i_1, i_2, i_3, i_4) \in C(I_M)} \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}}{\sqrt{\sum_{(i_1, i_2, i_3, i_4) \in C(I_N)} \text{Var}(\mathbf{W}_{i_1 i_2}) \text{Var}(\mathbf{W}_{i_2 i_3}) \text{Var}(\mathbf{W}_{i_3 i_4}) \text{Var}(\mathbf{W}_{i_4 i_1})}}.$$

It is easy to see that  $\mathbb{E}[S_{N,M} | \mathcal{F}_{N,M-1}] = S_{N,M-1}$ . Hence,  $\{Y_{N,M}\}_{M=1}^N$  is a martingale difference sequence relative to the filtration  $\{\mathcal{F}_{N,M}\}_{M=1}^N$ , and  $S_{N,N} = \sum_{M=1}^N Y_{N,M}$ . To show  $S_{N,N} \rightarrow \mathbb{N}(0, 1)$  as  $N \rightarrow \infty$ , we apply the martingale central limit theorem and check:

- (a)  $\sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) \rightarrow 1$  in probability

(b)  $\sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 \mathbf{1}_{\{|Y_{N,M}|>\epsilon\}} | \mathcal{F}_{N,M-1}) \rightarrow 0$ , in probability for any  $\epsilon > 0$ .

Note that once we have checked that both conditions (a) and (b) are satisfied, then by the martingale central limit theorem,  $S_{N,N} \rightarrow \mathbb{N}(0, 1)$ . Hence we have proved Lemma B.1.

It remains to check (a)-(b). For preparation, we first derive an alternative expression of  $\mathbb{E}(Y_{N,M} | \mathcal{F}_{N,M-1})$ . By definition,

$$Y_{N,M} = \frac{1}{\sqrt{D_N}} \sum_{(i_1, i_2, i_3, i_4) \in C(I_M) \setminus C(I_{M-1})} \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1},$$

where  $D_N = \sum_{(i_1, i_2, i_3, i_4) \in C(I_N)} \text{Var}(\mathbf{W}_{i_1 i_2}) \text{Var}(\mathbf{W}_{i_2 i_3}) \text{Var}(\mathbf{W}_{i_3 i_4}) \text{Var}(\mathbf{W}_{i_4 i_1})$ , and the summation is over all 4-cycles in  $C(I_M) \setminus C(I_{M-1})$ . Note that a cycle in  $C(I_M) \setminus C(I_{M-1})$  has to include the node  $M$ . Hence, we can use the following way to get all such cycles: First, select 2 indices  $(i, j)$  from  $\{1, 2, \dots, M-1\}$  and use them as the two neighboring nodes of  $M$ ; second, select an index  $k \in \{1, 2, \dots, M-1\} \setminus \{i, j\}$  as the last node in the cycle. This allows us to write

$$Y_{N,M} = \frac{1}{\sqrt{D_N}} \sum_{1 \leq i < j \leq M-1} \mathbf{W}_{Mi} \mathbf{W}_{Mj} \Gamma_{M-1,ij},$$

where

$$\Gamma_{M-1,ij} = \sum_{1 \leq k \leq M-1, k \notin \{i, j\}} \mathbf{W}_{ki} \mathbf{W}_{kj}.$$

Conditioning on  $\mathcal{F}_{N,M-1}$ ,  $\{\mathbf{W}_{Mi} \mathbf{W}_{Mj}\}_{1 \leq i < j \leq M-1}$  are mutually uncorrelated and  $\Gamma_{M-1,ij}$  is a constant. Hence, it follows that

$$\mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) = \frac{1}{D_N} \sum_{1 \leq i < j \leq M-1} \Gamma_{M-1,ij}^2 \text{Var}(\mathbf{W}_{Mi} \mathbf{W}_{Mj}) = \frac{1}{D_N} \sum_{1 \leq i < j \leq M-1} \Gamma_{M-1,ij}^2 \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}).$$

We now check (a). It suffices to show that

- (c)  $\mathbb{E} \left[ \sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) \right] = 1$
- (d)  $\text{Var} \left( \sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) \right) \rightarrow 0$ .

(Then (a) follows by Chebyshev inequality)

Consider (c). The terms in  $\Gamma_{M-1,ij}$  are unconditionally mutually uncorrelated. As a result,

$$\mathbb{E}[\Gamma_{M-1,ij}^2] = \sum_{k < M, k \notin \{i, j\}} \text{Var}(\mathbf{W}_{ki}) \text{Var}(\mathbf{W}_{kj}).$$

It follows that

$$\begin{aligned} \mathbb{E} \left[ \sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) \right] &= \frac{1}{D_N} \sum_{M=1}^N \sum_{1 \leq i < j \leq M-1} \sum_{1 \leq k \leq M-1, k \notin \{i, j\}} \text{Var}(\mathbf{W}_{ki}) \text{Var}(\mathbf{W}_{kj}) \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}) \\ &= \frac{1}{D_N} \sum_{(M,i,j,k) \in CC(I_N)} \text{Var}(\mathbf{W}_{ki}) \text{Var}(\mathbf{W}_{kj}) \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}) \\ &= 1. \end{aligned}$$

This proves (c).

Consider (d). We first decompose the random variable  $\sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1})$  into the sum of two parts, and then calculate its variance. We have

$$\Gamma_{M-1,ij}^2 = \sum_k \mathbf{W}_{ki}^2 \mathbf{W}_{kj}^2 + \sum_{k \neq \ell} \mathbf{W}_{ki} \mathbf{W}_{kj} \mathbf{W}_{\ell i} \mathbf{W}_{\ell j},$$

where  $k$  and  $\ell$  range in  $\{1, 2, \dots, M-1\} \setminus \{i, j\}$ . Now we can have a decomposition

$$\sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) = I_a + I_b,$$

where

$$I_a = \frac{1}{D_N} \sum_{M=1}^N \sum_{i < j \leq M-1} \sum_{k \leq M-1} \mathbf{W}_{ki}^2 \mathbf{W}_{kj}^2 \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}),$$

and

$$I_b = \frac{1}{D_N} \sum_{M=1}^N \sum_{i < j \leq M-1} \sum_{k, l \leq M-1} \sum_{k, l \notin \{i, j\}} \mathbf{W}_{ki} \mathbf{W}_{kj} \mathbf{W}_{\ell i} \mathbf{W}_{\ell j} \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}).$$

Then,

$$\text{Var} \left( \sum_{M=1}^N \mathbb{E}(Y_{N,M}^2 | \mathcal{F}_{N,M-1}) \right) = \text{Var}(I_a) + \text{Var}(I_b) + 2\text{Cov}(I_a, I_b) \leq (\sqrt{\text{Var}(I_a)} + \sqrt{\text{Var}(I_b)})^2.$$

It suffices to show that both  $\text{Var}(I_a) \rightarrow 0$  and  $\text{Var}(I_b) \rightarrow 0$ .

Consider the variance of  $I_a$ . In the sum of  $I_a$ , all 4-cycles  $(k, i, M, j)$  involved are selected in this way: We first select  $M$ , then select a pair  $(i, j)$  from  $\{1, 2, \dots, M-1\}$  and connect both  $i$  and  $j$  to  $M$ , and finally select  $k$  to close the cycle. In fact, these 4-cycles can be selected in an alternative way: First, select a V-shape  $(i, k, j)$  with  $k$  being the middle point. Second, select  $M > \max(i, k, j)$  to make the V-shape a cycle. Hence, we can rewrite

$$I_a = \frac{1}{D_N} \sum_{k=1}^N \sum_{1 \leq i < j \leq N} \sum_{i \neq k, j \neq k} \mathbf{W}_{ki}^2 \mathbf{W}_{kj}^2 \sum_{M > \max\{i, j, k\}} \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}).$$

$$(b_{kij} := \sum_{M > \max\{i, j, k\}} \text{Var}(\mathbf{W}_{Mi}) \text{Var}(\mathbf{W}_{Mj}))$$

We now fix  $k$  and calculate the covariance between  $\mathbf{W}_{ki}^2 \mathbf{W}_{kj}^2$  and  $\mathbf{W}_{ki'} \mathbf{W}_{kj'}$  for  $(i, j) \neq (i', j')$ . There are three cases.

Case (i):  $(i, j) = (i', j')$ . In this case,  $\text{Var}(\mathbf{W}_{ki}^2 \mathbf{W}_{kj}^2) \leq \mathbb{E}[\mathbf{W}_{ki}^4 \mathbf{W}_{kj}^4] = \mathbb{E}[\mathbf{W}_{ki}^4] \mathbb{E}[\mathbf{W}_{kj}^4] \leq C\sigma^8$ .

Case (ii):  $i = i'$  but  $j \neq j'$ . In this case, we have  $\text{Cov}(\mathbf{W}_{ki}^2 \mathbf{W}_{kj}^2, \mathbf{W}_{ki}^2 \mathbf{W}_{kj'}^2) = \text{Var}(\mathbf{W}_{ki}^2) \mathbb{E}[\mathbf{W}_{kj}^2 \mathbf{W}_{kj'}^2] \leq C\sigma^4 \cdot \sigma^2 \cdot \sigma^2$ .

Case (iii):  $i \neq i'$  and  $j \neq j'$ . The two terms are independent, and their covariance is zero. Combining the above gives

$$\text{Var}(I_a) \leq \frac{N}{D_N^2} \sum_{k=1}^N \left( \sum_{1 \leq i < j \leq N} \sum_{k \neq i, k \neq j} b_{kij}^2 \cdot C\sigma^8 + \sum_{i, j, j' \in \{1, \dots, N\} \setminus k} \sum_{i, j, j' \text{ are distinct}} b_{kij} b_{kij'} \cdot \sigma^8 \right).$$

(Here we use the inequality  $\text{Var}(\sum_{i=1}^N Y_i) \leq N \sum_{i=1}^N \text{Var}(Y_i)$ )

We now bound the right hand side.  $\text{Var}(\mathbf{W}_{ij}) \leq \sigma^2$ . Hence,  $b_{kij} \leq \sum_{M>k} \sigma^4 \leq N\sigma^4$ . As a result,

$$\text{Var}(I_a) \leq \frac{N}{D_N^2} \left( \sum_{k,i,j} 256N^2\sigma^{16} + \sum_{k,i,j,j'} 16N^2\sigma^{16} \right) \leq \frac{CN^7\sigma^{16}}{D_N^2}.$$

Moreover, since  $D_N \geq \tau^8 N(N-1)(N-2)(N-3)$ . As a result, we have

$$\text{Var}(I_a) \leq \frac{C\sigma^{16}N^7}{\tau^{16}N^8} = o(1).$$

Consider the variance of  $I_b$ . Rewrite

$$I_b = \frac{1}{D_N} \sum_{k,j,l,i \text{ dist}} c_{klij} G_{klij}.$$

where  $G_{klij} := \mathbf{W}_{ki}\mathbf{W}_{kj}\mathbf{W}_{li}\mathbf{W}_{lj}$ ,  $c_{klij} := \sum_{M>\max\{k,l,i,j\}} \text{Var}(\mathbf{W}_{Mi})\text{Var}(\mathbf{W}_{Mj})$ .

Since  $I_b$  has a mean zero,  $\text{Var}(I_b) = \mathbb{E}(I_b^2)$ . Additionally, for 2 cycles  $(k, \ell, i, j)$  and  $(k', \ell', i', j')$ , only when they are exactly equal, we have  $\mathbb{E}[G_{klij}G_{k'\ell'i'j'}] \neq 0$ . As a result,

$$\text{Var}(I_b) = \frac{1}{D_N^2} \sum_{k,\ell,i,j \text{ are distinct}} c_{klij}^2 \mathbb{E}[G_{klij}^2] = \frac{1}{D_N^2} \sum_{k,\ell,i,j \text{ are distinct}} c_{klij}^2 \text{Var}(\mathbf{W}_{ki})\text{Var}(\mathbf{W}_{kj})\text{Var}(\mathbf{W}_{li})\text{Var}(\mathbf{W}_{lj}).$$

Similarly to how we get the bound for  $b_{ijk}$ , we can derive that  $c_{klij} \leq N\sigma^4$ . Hence,

$$\text{Var}(I_b) \leq \frac{1}{\tau^{16}N^8} N^4 \cdot N^2\sigma^{16} = o(1).$$

As a result,

$$\sqrt{\text{Var}(I_b)} = o(1).$$

Thus we have proved (a).

We now check (b). By the Cauchy-Schwarz inequality and the Chebyshev's inequality,

$$\begin{aligned} \sum_{M=1}^N \mathbb{E} \left[ Y_{N,M}^2 \mathbf{1}_{\{Y_{N,M}>\epsilon\}} \mid \mathcal{F}_{N,M-1} \right] &\leq \sum_{M=1}^N \sqrt{\mathbb{E} \left[ Y_{N,M}^4 \mid \mathcal{F}_{N,M-1} \right]} \sqrt{\mathbb{P}(Y_{N,M} \geq \epsilon \mid \mathcal{F}_{N,M-1})} \\ &\leq \epsilon^{-2} \sum_{M=1}^N \mathbb{E} \left[ Y_{N,M}^4 \mid \mathcal{F}_{N,M-1} \right]. \end{aligned}$$

Therefore, it suffices to show that the right-hand side converges to zero in probability. Then, it suffices to show that its  $L^1$ -norm converges to zero. Since the right-hand side is a nonnegative random variable, we only need to prove that its expectation converges to zero, i.e.,

$$\mathbb{E} \left[ \sum_{M=1}^N Y_{N,M}^4 \right] = o(1).$$

We have

$$\mathbb{E}[Y_{N,M}^4] = \frac{1}{D_N^2} \left( \sum_{1 \leq i < j \leq M-1} \mathbb{E}[\mathbf{W}_{Mi}^4 \mathbf{W}_{Mj}^4] \mathbb{E}[\Gamma_{M-1,ij}^4] + 4 \sum_{(i,j) \neq (i',j')} \mathbb{E}[\Gamma_{M-1,ij}^2 \Gamma_{M-1,i'j'}^2] \mathbb{E}[\mathbf{W}_{Mi}^2 \mathbf{W}_{Mj}^2 \mathbf{W}_{Mi'}^2] \right)$$

$$\mathbf{W}_{Mj'}^2])$$

since  $\mathbb{E}[\mathbf{W}_{Mi}\mathbf{W}_{Mj}\mathbf{W}_{Mi'}\mathbf{W}_{Mj'}\mathbf{W}_{Mi''}\mathbf{W}_{Mj''}\mathbf{W}_{Mi'''}\mathbf{W}_{Mj'''}]=0$  if any of  $i, j, i', j', i'', j'', i''', j'''$  appear only once. Note that if  $(i, j) \neq (i', j')$ ,  $\mathbb{E}[\mathbf{W}_{k_1i}\mathbf{W}_{k_1j}\mathbf{W}_{k_2i}\mathbf{W}_{k_2j}\mathbf{W}_{k_3i'}\mathbf{W}_{k_3j'}\mathbf{W}_{k_4i'}\mathbf{W}_{k_4j'}]=0$  unless  $k_1 = k_2, k_3 = k_4$ , then we have

$$\begin{aligned}\mathbb{E}[\Gamma_{M-1,ij}^2\Gamma_{M-1,i'j'}^2] &\leq \sum_{1 \leq k_1, k_2, k_3, k_4 \leq M-1} \sum_{k_1, k_2 \notin \{i, j\}, k_3, k_4 \notin \{i', j'\}} \mathbb{E}[\mathbf{W}_{k_1i}\mathbf{W}_{k_1j}\mathbf{W}_{k_2i}\mathbf{W}_{k_2j}\mathbf{W}_{k_3i'}\mathbf{W}_{k_3j'}\mathbf{W}_{k_4i'}\mathbf{W}_{k_4j'}] \\ &= \sum_{k_1=k_2, k_3=k_4} \mathbb{E}[G_{k_1k_2ij}G_{k_3k_4i'j'}] \\ &\leq CN^2\sigma^8\end{aligned}$$

Combining this with the fact that  $\mathbb{E}[\Gamma_{M-1,ij}^4] \leq CN^4\sigma^8$  (each term is smaller than  $C\sigma^8$ ), we deduce that

$$\mathbb{E}[Y_{N,M}^4] \leq \frac{1}{\tau^{16}N^8} \cdot CN^6\sigma^{16} = \frac{\sigma^{16}}{\tau^{16}N^2}$$

As a result,

$$\sum_{M=1}^N \mathbb{E}[Y_{N,M}^4] \leq \frac{\sigma^{16}}{\tau^{16}N} = o(1).$$

This gives (b) follows.

**Remark 6.** We largely follow the proof presented in Jin et al. (2018). However, their proof for assertion (b) might contain a typo due to a miscalculation of  $Y_{N,M}^4$ , which is denoted as  $X_{N,M}^4$  there. To be precise, according to Equation (12) in the Supplement of Jin et al. (2018), they define  $X_{n,m} = \frac{1}{\sqrt{M_n}} \sum_{1 \leq i < j \leq m-1} W_{mi}W_{mj}Y_{(m-1)ij}$ , where  $M_n$  is a constant. However, in the subsequent derivation of  $\mathbb{E}[X_{n,m}^4 | \mathcal{F}_{n,m-1}]$ , the denominator is presented as  $M_n^4$ , which I was unable to obtain, at least, based on the earlier result. In any case, this appears to be a minor typographical error that can be easily corrected.

To address this issue, we have modified their proof to ensure correctness.

### S.3 Proof of Lemma S.2

The proof is combined with the proof of Lemma C.3, see below.

### S.4 Proof of Theorem 3

Recall that  $Z$  is the true community label matrix. Fix  $1 \leq m < K$ . Let  $\{\mathcal{G}_m\}$  be the class of  $N \times m$  matrices  $\mathbf{Z}^{(0)}$ , where each  $\mathbf{Z}^{(0)}$  is formed as follows: let  $\{1, 2, \dots, K\} = S_1 \cup S_2 \dots \cup S_m$  be a partition, column  $\ell$  of  $Z^{(0)}$  is the sum of all columns of  $\mathbf{Z}$  in  $S_\ell$ ,  $1 \leq \ell \leq m$ . Let  $L^{(0)}$  be the  $K \times m$  matrix of 0 and 1 where

$$L^{(0)}(k, \ell) = 1 \text{ if and only if } k \text{ in } S_\ell, 1 \leq k \leq K, 1 \leq \ell \leq m.$$

Therefore, for each  $\mathbf{Z}^{(0)} \in \mathcal{G}_m$ , we can find an  $L^{(0)}$  such that  $\mathbf{Z}^{(0)} = \mathbf{Z}L^{(0)}$ .

Now we can construct  $\widehat{\mathbf{P}}^{(m,0)}$  based on  $\mathbf{Z}^{(0)}$  and introduce  $\widehat{\mathbf{P}}^{(m,0)} = \mathbf{Z}^{(0)} ((\mathbf{Z}^{(0)})^\top \mathbf{Z}^{(0)})^{-1} (\mathbf{Z}^{(0)})^\top \mathbf{X}$ ,  $Q_N^{(m,0)} = \sum_{i_1, i_2, i_3, i_4 \text{ (dist)}} (\mathcal{S}(\mathbf{X})_{i_1, i_2} - (\mathcal{S}(\widehat{\mathbf{P}}^{(m,0)})_{i_1, i_2})) (\mathcal{S}(\mathbf{X})_{i_2, i_3} - (\mathcal{S}(\widehat{\mathbf{P}}^{(m,0)})_{i_2, i_3})) (\mathcal{S}(\mathbf{X})_{i_3, i_4} - (\mathcal{S}(\widehat{\mathbf{P}}^{(m,0)})_{i_3, i_4})) (\mathcal{S}(\mathbf{X})_{i_4, i_1} - (\mathcal{S}(\widehat{\mathbf{P}}^{(m,0)})_{i_4, i_1}))$ ,

and  $\phi_N^{(m,0)} = Q_N^{(m,0)}/\sqrt{C_N}$ . These are the proxies of  $\widehat{\mathbf{P}}^{(m)}$ ,  $Q_N^{(m)}$  and  $\phi_N^{(m)}$ , respectively, where  $\widehat{\mathbf{Z}}^{(m)}$  is now frozen at  $\mathbf{Z}^{(0)}$ .

Now we define a non-stochastic counterpart of  $\widehat{\mathbf{P}}^{(m,0)}$  as follows. Let  $\mathbf{P}^{(m,0)}$  be constructed similarly to  $\widehat{\mathbf{P}}^{(m,0)}$ , except that  $\mathbf{X}$  is replaced with  $\mathbf{P}$ . Similarly, we can define the following proxy of  $Q_N^{(m,0)}$ .

$$\widetilde{Q}_N^{(m,0)} = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} (\mathcal{S}(\mathbf{X})_{i_1, i_2} - (\mathcal{S}(\mathbf{P}^{(m,0)})_{i_1, i_2})(\mathcal{S}(\mathbf{X})_{i_2, i_3} - (\mathcal{S}(\mathbf{P}^{(m,0)})_{i_2, i_3})(\mathcal{S}(\mathbf{X})_{i_3, i_4} - (\mathcal{S}(\mathbf{P}^{(m,0)})_{i_3, i_4})) \\ (\mathcal{S}(\mathbf{X})_{i_4, i_1} - (\mathcal{S}(\mathbf{P}^{(m,0)})_{i_4, i_1}))$$

Introduce  $\widetilde{\mathbf{P}}^{(m,0)} = \mathbf{P} - \mathbf{P}^{(m,0)}$ , thus we can rewrite  $\widetilde{Q}_N^{(m,0)}$  as

$$\widetilde{Q}_N^{(m,0)} = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} (\mathbf{W}_{i_1 i_2} + \mathcal{S}(\widetilde{\mathbf{P}}^{(m,0)})_{i_1 i_2})(\mathbf{W}_{i_2 i_3} + \mathcal{S}(\widetilde{\mathbf{P}}^{(m,0)})_{i_2 i_3})(\mathbf{W}_{i_3 i_4} + \mathcal{S}(\widetilde{\mathbf{P}}^{(m,0)})_{i_3 i_4}) \\ (\mathbf{W}_{i_4 i_1} + \mathcal{S}(\widetilde{\mathbf{P}}^{(m,0)})_{i_4 i_1}).$$

Let  $\tilde{\sigma}_k$  be the  $k$ -th largest (in magnitude) singular value of  $\widetilde{\mathbf{P}}^{(m,0)}$  and recall that  $\sigma_k$  is the  $k$ -th largest (in magnitude) singular value of  $\mathbf{P}$ . We have following lemmas.

**Lemma S.3.** *For each  $1 \leq m \leq K$ ,  $\text{tr}(\mathcal{S}(\widetilde{\mathbf{P}}^{(m,0)})^4) \geq C\sigma^4 N^4$ .*

**Lemma S.4.** *For  $1 \leq m < K$ ,*

$$\mathbb{E}[\widetilde{Q}_N^{(m,0)}] = \text{tr}(\mathcal{S}(\widetilde{\mathbf{P}}^{(m,0)})^4) + o(N^4), \quad \text{Var}(\widetilde{Q}_N^{(m,0)}) \leq C(N^6\sigma^2 + N^5\sigma^4 + N^4\sigma^8).$$

**Lemma S.5.** *For  $1 \leq m < K$ ,*

$$|\mathbb{E}[Q_N^{(m,0)} - \widetilde{Q}_N^{(m,0)}]| = O(\sigma^4 p^2), \quad \text{Var}(Q_N^{(m,0)} - \widetilde{Q}_N^{(m,0)}) = o(N^8).$$

For notation simplicity, we write  $\widetilde{\mathbf{P}}^{(m,0)} = \widetilde{\mathbf{P}}$ .

We now prove Theorem 3. Note that by Theorem 2, the second item of Theorem 3 follows once the first item is proved. Therefore we only consider the first item, where it is sufficient to show that for all  $1 < m < K$ ,

$$\phi_N^{(m)} \rightarrow \infty, \quad \text{in probability.}$$

By Theorem 1, there exists an event  $A_n$  with  $\mathbb{P}(A_n^c) \leq Cn^{-5}$  as  $n \rightarrow \infty$ , such that on event  $A_n$  we have  $\widehat{\mathbf{Z}}^{(m)} \in \mathcal{G}_m$ . This further indicates that on event  $A_n$  we have

$$\phi_N^{(m)} \geq \min_{Z^{(0)} \in \mathcal{G}_m} \phi_N^{(m,0)}$$

Then further notice that the cardinality of  $\mathcal{G}_m$  are  $m^K$ , which is of constant order as long as  $K$  is constant. Therefore to prove  $\phi_N^{(m)} \rightarrow \infty$  in probability, it suffices to show that for any fixed  $\mathbf{Z}^{(0)} \in \mathcal{G}_m$ .

$$\phi_N^{(m,0)} \rightarrow \infty, \quad \text{in probability.} \tag{S.1}$$

By Lemma S.3-S.5,

$$\mathbb{E}\left[\frac{Q_N^{(m,0)}}{\sqrt{C_N}}\right] \geq CN^{4-\frac{\beta}{2}} \cdot [1 + o(1)] \rightarrow \infty, \quad \text{Var}\left(\frac{Q_N^{(m,0)}}{\sqrt{C_N}}\right) \leq CN^{7-\beta}.$$

Therefore, by Chebyshev's inequality, for any constant  $M > 0$ ,

$$\mathbb{P}\left(\frac{Q_N^{(m,0)}}{\sqrt{C_N}} < M\right) \leq (\mathbb{E}\left[\frac{Q_N^{(m,0)}}{\sqrt{C_N}}\right] - M)^{-2} \text{Var}\left(\frac{Q_N^{(m,0)}}{\sqrt{C_N}}\right) \leq C \left[ \frac{N^{7-\beta}}{(N^{4-\frac{\beta}{2}}[1 + o(1)] - M)^2} \right] \rightarrow 0,$$

Hence we conclude the proof of Theorem 3.

**Remark 7.** Based on Assumptions 1 and 4,  $\sigma$  is a constant and each component of the cluster centers is bounded by  $C_P$ . Combining this with the Non-Splitting Property, it follows that  $\hat{\sigma}$ , obtained from Algorithm 2, is also bounded with high probability. Since the proofs of Theorems 2 and 3 are based on the order of  $N$ , the boundedness of  $\hat{\sigma}$  ensures that replacing  $\sigma$  with  $\hat{\sigma}$  does not affect the validity of our results. Consequently, the proofs for Algorithm 1 also apply to Algorithm 2.

## S.5 Proof of Lemma S.3

By definition, it is easy to see

$$((\widehat{\mathbf{Z}}^{(m)})^\top \widehat{\mathbf{Z}}^{(m)})^{-1} (\widehat{\mathbf{Z}}^{(m)})^\top \mathbf{P} = \begin{pmatrix} a_{11}\boldsymbol{\theta}_1^* + \cdots + a_{1K}\boldsymbol{\theta}_K^* \\ \vdots \\ a_{m1}\boldsymbol{\theta}_1^* + \cdots + a_{mK}\boldsymbol{\theta}_K^* \end{pmatrix},$$

where  $a_{11}, \dots, a_{mK} \in [0, 1]$  satisfy  $\sum_{j=1}^K a_{ij} = 1$ ,  $i = 1, \dots, m$ .

It follows that the 2-norm of each row in  $\tilde{\mathbf{P}} = \mathbf{P} - \mathbf{P}^{(m)}$  is the Euclidean distance from one of  $\boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_K^*$  to one of  $a_{11}\boldsymbol{\theta}_1^* + \cdots + a_{1K}\boldsymbol{\theta}_K^*, \dots, a_{m1}\boldsymbol{\theta}_1^* + \cdots + a_{mK}\boldsymbol{\theta}_K^*$ .

Recall that  $\mathbf{z}_1, \dots, \mathbf{z}_n$  are true row cluster assignments, we introduce  $\widehat{\mathbf{z}}_1^{(m)}, \dots, \widehat{\mathbf{z}}_n^{(m)}$  are pseudo row cluster assignments and pseudo row clusters  $C_1^{(m)}, \dots, C_m^{(m)}$  corresponding to  $\widehat{\mathbf{Z}}^{(m)}$ . To this end, we can rewrite  $\tilde{\mathbf{P}}$  as

$$\begin{pmatrix} \boldsymbol{\theta}_{\mathbf{z}_1}^* - \sum_{j=1}^K a_{\widehat{\mathbf{z}}_1^{(m)} j} \boldsymbol{\theta}_j^* \\ \vdots \\ \boldsymbol{\theta}_{\mathbf{z}_n}^* - \sum_{j=1}^K a_{\widehat{\mathbf{z}}_n^{(m)} j} \boldsymbol{\theta}_j^* \end{pmatrix}.$$

For any  $i = 1, \dots, K$ , using the pigeonhole principle and

$$|C_i| = |C_i \cap C_1^{(m)}| + \cdots + |C_i \cap C_m^{(m)}|,$$

we can deduce there exists  $t_i \in \{1, \dots, m\}$  such that

$$|C_i \cap C_{t_i}^{(m)}| \geq \frac{|C_i|}{m} \geq \frac{\alpha_0}{K} n.$$

As a result,  $\boldsymbol{\theta}_1^* - \sum_{j=1}^K a_{t_1 j} \boldsymbol{\theta}_j^*, \dots, \boldsymbol{\theta}_K^* - \sum_{j=1}^K a_{t_K j} \boldsymbol{\theta}_j^*$  appear at least  $\frac{\alpha_0}{K} n$  times across all the rows of  $\tilde{\mathbf{P}}$ .

Since  $t_1, \dots, t_K \in \{1, \dots, m\}$ , using pigeonhole principle again, we deduce that there exist  $u \neq v$  such that  $t_u = t_v$ . Therefore,  $\boldsymbol{\theta}_u^* - \sum_{j=1}^K a_{t_u j} \boldsymbol{\theta}_j^*$  and  $\boldsymbol{\theta}_v^* - \sum_{j=1}^K a_{t_v j} \boldsymbol{\theta}_j^*$  both appear at least  $\frac{\alpha_0}{K} n$  times across all the rows of  $\tilde{\mathbf{P}}$ .

By triangle inequality, we know

$$\max \left\{ \|\boldsymbol{\theta}_u^* - \sum_{j=1}^K a_{t_{uj}} \boldsymbol{\theta}_j^*\|_2, \|\boldsymbol{\theta}_v^* - \sum_{j=1}^K a_{t_{uj}} \boldsymbol{\theta}_j^*\|_2 \right\} \geq \frac{1}{2} \|\boldsymbol{\theta}_u^* - \boldsymbol{\theta}_v^*\| \geq \frac{1}{2} \Delta.$$

Without loss of generality, assume  $\|\boldsymbol{\theta}_u^* - \sum_{j=1}^K a_{t_{uj}} \boldsymbol{\theta}_j^*\|_2$  is the larger one. Since it appears at least  $\frac{\alpha_0}{K} n$  times across all the rows of  $\tilde{\mathbf{P}}$ , we can find a submatrix of  $\tilde{\mathbf{P}}$  consisting of  $\frac{\alpha_0}{K} n \|\boldsymbol{\theta}_u^* - \sum_{j=1}^K a_{t_{uj}} \boldsymbol{\theta}_j^*\|_2$ 's vertically stacked together. It is easy to see the 2-norm of this submatrix is larger than  $\frac{1}{2} \sqrt{\frac{\alpha_0}{K} n} \Delta$ . We conclude that

$$\|\tilde{\mathbf{P}}\|_2 = \omega(\sqrt{n}\Delta) = \omega(n\sigma).$$

As a result,  $\text{tr}(\mathcal{S}(\tilde{\mathbf{P}}^{(m,0)})^4) \geq C\sigma^4 N^4$ .

**Remark 8.** If we further assume  $\kappa(\boldsymbol{\Theta}) = O(1)$ , this lemma follows directly from the facts that  $\text{rank}(\widehat{\mathbf{Z}}^{(m)}) = m$  and  $\sigma_K = \omega(\sqrt{N})$

## S.6 Proof of Lemma S.4

Given an  $N \times N$  symmetric matrix  $\mathbf{T}$ , we define a random variable:

$$\mathcal{Q}_W(\mathbf{T}) = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} (\mathbf{W}_{i_1 i_2} + \mathbf{T}_{i_1 i_2})(\mathbf{W}_{i_2 i_3} + \mathbf{T}_{i_2 i_3})(\mathbf{W}_{i_3 i_4} + \mathbf{T}_{i_3 i_4})(\mathbf{W}_{i_4 i_1} + \mathbf{T}_{i_4 i_1}).$$

Then,  $\tilde{Q}_N^{(m,0)}$  is a special case with  $\mathbf{T} = \mathcal{S}(\tilde{\mathbf{P}})$ . We aim to study the general form of  $\mathcal{Q}_W(\mathbf{T})$  and prove the following lemma:

**Lemma S.6.** As  $N \rightarrow \infty$ , suppose there is a constant  $C > 0$  such that  $|\mathbf{T}_{ij}| \leq C$  for all  $1 \leq i, j \leq N$ . Then,  $\mathbb{E}[\mathcal{Q}_W(\mathbf{T})] = \text{tr}(\mathbf{T}^4) + o(N^4)$  and  $\text{Var}(\mathcal{Q}_W(\mathbf{T})) \leq CN^6$ .

We now set  $\mathbf{T} = \mathcal{S}(\tilde{\mathbf{P}}^{(m,0)})$  and verify the conditions of Lemma S.6. By Assumption 7,  $\mathcal{S}(\tilde{\mathbf{P}}^{(m)}) \leq 2C_P$  and hence we can apply this lemma. The claim follows immediately.

It remains to show Lemma S.6. We write  $\mathcal{Q}_W(\mathbf{T})$  as the sum of  $2^4 = 16$  post-expansion sums. Each post-expansion sum takes a form

$$X = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} a_{i_1 i_2} b_{i_2 i_3} c_{i_3 i_4} d_{i_4 i_1},$$

where each of  $a_{ij}, b_{ij}, c_{ij}, d_{ij}$  may take value in  $\{\mathbf{W}_{ij}, \mathbf{T}_{ij}\}$ . Then,  $\mathbb{E}[X]$  is equal to the sum of means of these post-expansion sums, and  $\text{Var}(X)$  is bounded by a constant times the sum of variances of these post-expansion sums. It suffices to study the means and variances of these post-expansion sums.

We divide 16 post-expansion sums into 6 common types and compute the mean and variance of each type.

The calculation of mean is easy since  $\mathbf{W}_{ij}$  and  $\mathbf{W}_{i'j'}$  are independent if  $(i, j)$  and  $(i', j')$  are distinct. Besides, we have  $\mathbb{E}[X_6] = \text{tr}(\mathbf{T}^4) - \Delta$  and  $|\mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{T}_{i_4 i_1}| \leq C$ , then it follows that  $\mathbb{E}[X_6] = \text{tr}(\mathbf{T}^4) + C \sum_{(i_1, i_2, i_3, i_4) \text{non-dist}} 1 = \text{tr}(\mathbf{T}^4) + O(N^3)$ .

Now we turn to the calculation of variance.

We already see  $\text{Var}(X_1) \leq CN^4 \sigma^8$  in the proof of Theorem 4.

Now we introduce three terms below.

| Type | # | Examples  | Mean                      | Variance          |
|------|---|---|---------------------------|-------------------|
| I    | 1 | $X_1 = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$ | 0                         | $O(N^4 \sigma^8)$ |
| II   | 4 | $X_2 = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{T}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$ | 0                         | $O(N^4 \sigma^6)$ |
| IIIa | 4 | $X_3 = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$ | 0                         | $O(N^5 \sigma^4)$ |
| IIIb | 2 | $X_4 = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{T}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$ | 0                         | $O(N^5 \sigma^4)$ |
| IV   | 4 | $X_5 = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$ | 0                         | $O(N^6 \sigma^2)$ |
| V    | 1 | $X_6 = \sum_{i_1, i_2, i_3, i_4 (\text{dist})} \mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{T}_{i_4 i_1}$ | $\text{tr}(T^4) + o(N^4)$ | 0                 |

Table S.1: The post-expansion sums of  $\mathcal{Q}_W(T)$  have 6 different types. We present the mean and variance of each type.

$$\begin{aligned} \chi_{i_1, i_2, i_3, i_4}^{(1)} &= \mathbf{W}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{T}_{i_4 i_1} + \mathbf{T}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{T}_{i_4 i_1} \\ &\quad + \mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{T}_{i_4 i_1} + \mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{W}_{i_4 i_1}, \end{aligned}$$

$$\begin{aligned} \chi_{i_1, i_2, i_3, i_4}^{(2)} &= \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{T}_{i_4 i_1} + \mathbf{W}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{T}_{i_4 i_1} + \mathbf{W}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \\ &\quad + \mathbf{T}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{T}_{i_4 i_1} + \mathbf{T}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{T}_{i_4 i_1} + \mathbf{T}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}, \end{aligned}$$

$$\chi_{i_1, i_2, i_3, i_4}^{(3)} = \mathbf{T}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} + \mathbf{W}_{i_1 i_2} \mathbf{T}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} + \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{T}_{i_3 i_4} \mathbf{W}_{i_4 i_1} + \mathbf{W}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{T}_{i_4 i_1}.$$

Note that the four terms in  $\chi_{i_1, i_2, i_3, i_4}^{(1)}$  are independent of each other. Hence,  $\text{Var}(\chi_{i_1, i_2, i_3, i_4}^{(1)}) \leq C\sigma^2$  and  $\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \text{Var}(\chi_{i_1, i_2, i_3, i_4}^{(1)}) \leq CN^4\sigma^2$ .

We then look at the covariance between  $\chi_{i_1, i_2, i_3, i_4}^{(1)}$  and  $\chi_{i'_1, i'_2, i'_3, i'_4}^{(1)}$ . Let  $(j, s, m, l)$  be any cycle on the four nodes  $\{i_1, i_2, i_3, i_4\}$ , and let  $(j', s', m', l')$  be any cycle on the four nodes  $\{i'_1, i'_2, i'_3, i'_4\}$ . As long as  $\{j, s\} \neq \{j', s'\}$ , the two terms  $\mathbf{W}_{js} \mathbf{T}_{sm} \mathbf{T}_{ml} \mathbf{T}_{lj}$  and  $\mathbf{W}_{j's'} \mathbf{T}_{s'm'} \mathbf{T}_{m'l'} \mathbf{T}_{l'j'}$  are independent, hence, their covariance is zero. Otherwise, the covariance is bounded by  $C\sigma^2$ . It follows that

$$\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \sum_{(i'_1, i'_2, i'_3, i'_4) \text{ dist}} \text{Cov}(\chi_{i_1, i_2, i_3, i_4}^{(1)}, \chi_{i'_1, i'_2, i'_3, i'_4}^{(1)}) \leq C \sum_{i_1, i_2, i_3, i_4, i'_3, i'_4} \sigma^2 \leq CN^6$$

Hence, the above imply  $\text{Var}(\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \chi_{i_1, i_2, i_3, i_4}^{(1)}) \leq CN^6$ .

We can consider other terms similarly. We have  $\mathbb{E}[(\chi_{i_1, i_2, i_3, i_4}^{(2)})] = 0$ ,  $\mathbb{E}[(\chi_{i_1, i_2, i_3, i_4}^{(3)})] = 0$ ,  $\text{Var}(\chi_{i_1, i_2, i_3, i_4}^{(2)}) \leq \mathbb{E}[(\chi_{i_1, i_2, i_3, i_4}^{(2)})^2] \leq 36\sigma^4$ ,  $\text{Var}(\chi_{i_1, i_2, i_3, i_4}^{(3)}) \leq \mathbb{E}[(\chi_{i_1, i_2, i_3, i_4}^{(3)})^2] \leq 16\sigma^6$ ,  $\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \text{Var}(\chi_{i_1, i_2, i_3, i_4}^{(2)}) \leq CN^4\sigma^4$  and  $\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \text{Var}(\chi_{i_1, i_2, i_3, i_4}^{(3)}) \leq CN^4\sigma^6$ . Additionally, to ensure  $\text{Cov}(\chi_{i_1, i_2, i_3, i_4}^{(2)}, \chi_{i'_1, i'_2, i'_3, i'_4}^{(2)})$  (resp.  $\text{Cov}(\chi_{i_1, i_2, i_3, i_4}^{(3)}, \chi_{i'_1, i'_2, i'_3, i'_4}^{(3)})$ ) is nonzero, we need  $\#\{(i_1, i_2, i_3, i_4) \cap (i'_1, i'_2, i'_3, i'_4)\} \geq 3$  (resp.  $\#\{(i_1, i_2, i_3, i_4) \cap (i'_1, i'_2, i'_3, i'_4)\} = 4$ ) and hence

$$\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \sum_{(i'_1, i'_2, i'_3, i'_4) \text{ dist}} \text{Cov}(\chi_{i_1, i_2, i_3, i_4}^{(2)}, \chi_{i'_1, i'_2, i'_3, i'_4}^{(2)}) \leq CN^5,$$

$$\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \sum_{(i'_1, i'_2, i'_3, i'_4) \text{ dist}} \text{Cov}(\chi_{i_1, i_2, i_3, i_4}^{(3)}, \chi_{i'_1, i'_2, i'_3, i'_4}^{(3)}) \leq CN^4,$$

which imply  $\text{Var}(\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \chi_{i_1, i_2, i_3, i_4}^{(2)}) \leq CN^5$  and  $\text{Var}(\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \chi_{i_1, i_2, i_3, i_4}^{(3)}) \leq CN^4$ .

Therefore, for  $\chi = \sum_{(i_1, i_2, i_3, i_4) \text{ dist}} (\chi_{i_1, i_2, i_3, i_4}^{(1)} + \chi_{i_1, i_2, i_3, i_4}^{(2)} + \chi_{i_1, i_2, i_3, i_4}^{(3)})$ , we have

$$\begin{aligned} \text{Var}(\chi) &\leq 3(\text{Var}(\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \chi_{i_1, i_2, i_3, i_4}^{(1)}) + \text{Var}(\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \chi_{i_1, i_2, i_3, i_4}^{(2)}) + \text{Var}(\sum_{(i_1, i_2, i_3, i_4) \text{ dist}} \chi_{i_1, i_2, i_3, i_4}^{(3)})) \\ &\leq CN^6 \end{aligned}$$

Consequently,

$$\text{Var}(Q_W(\mathbf{T})) \leq 2(\text{Var}(X_1) + \text{Var}(\chi)) \leq CN^6.$$

## S.7 Proof of Lemma S.5

Recall that our objective is to analyze the quantities  $|\mathbb{E}[Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}]|$  and  $\text{Var}(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)})$ . To this end, we first examine the expression  $Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}$ .

We introduce the notation  $M_{ijkl}(\mathbf{X}) = \mathbf{X}_{ij}\mathbf{X}_{jk}\mathbf{X}_{kl}\mathbf{X}_{\ell i}$  for any symmetric matrix  $\mathbf{X}$  and distinct indices  $(i, j, k, \ell)$ . Thus, we have

$$Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)} = \sum_{i_1, i_2, i_3, i_4 \text{ (dist)}} [M_{i_1 i_2 i_3 i_4}(\mathbf{X}) - M_{i_1 i_2 i_3 i_4}(\tilde{\mathbf{X}})],$$

where

$$\begin{cases} \mathbf{X}_{ij} = \mathcal{S}(\tilde{\mathbf{P}}^{(m,0)})_{ij} + \mathbf{W}_{ij} + \mathbf{D}_{ij}, \\ \tilde{\mathbf{X}}_{ij} = \mathcal{S}(\tilde{\mathbf{P}}^{(m,0)})_{ij} + \mathbf{W}_{ij}. \end{cases}$$

The matrix  $\mathbf{D}_{ij}$  is defined as

$$\mathbf{D}_{ij} = \begin{cases} 0, & \text{if } (i \leq n, j \leq n) \text{ and } (i > n, j > n), \\ \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_i}^{(m)}} \mathbf{W}_{kj}}{|C_{\tilde{\mathbf{z}}_i}^{(m)}|}, & \text{if } (i \leq n, j > n), \\ \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_j}^{(m)}} \mathbf{W}_{il}}{|C_{\tilde{\mathbf{z}}_j}^{(m)}|}, & \text{if } (i > n, j \leq n). \end{cases}$$

For simplicity, we shall omit the superscripts  $(m, 0)$  in  $(\tilde{\mathbf{P}}, \epsilon)$ . From the expressions for  $\mathbf{X}_{ij}$  and  $\tilde{\mathbf{X}}_{ij}$ , we observe that  $M_{i_1 i_2 i_3 i_4}(\mathbf{X}) - M_{i_1 i_2 i_3 i_4}(\tilde{\mathbf{X}})$  expands into  $3^4 - 2^4 = 65$  terms. Therefore,  $Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}$  consists of 65 post-expansion sums, each of the form

$$\sum_{(i_1, i_2, i_3, i_4) \text{ (dist)}} a_{i_1 i_2} b_{i_2 i_3} c_{i_3 i_4} d_{i_4 i_1}, \quad \text{where } a, b, c, d \in \{\mathcal{S}(\tilde{\mathbf{P}}), \mathbf{W}, \mathbf{D}\}.$$

To analyze  $|\mathbb{E}[Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}]|$  and  $\text{Var}(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)})$ , we apply the triangle inequality and Cauchy inequality, reducing the problem to evaluating the absolute mean and variance of each of these 65 post-expansion sums. In Table S.2, we categorize them into 15 distinct types, displaying their respective counts, absolute means and variances.

We shall proceed to verify each result in Table S.2 individually. Additionally, we will establish an upper bound for the second moment of each type, which will allow us to control the variance.

Actually, our primary objective is to prove that  $\text{Var}(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}) = o(N^8)$ . Accordingly, it suffices to demonstrate that  $\mathbb{E}[Y_i^2], \mathbb{E}[Z_j^2], \mathbb{E}[T_k^2], \mathbb{E}[F^2] = o(N^8)$  for  $i, j = 1, \dots, 6$  and  $k = 1, 2$ .

Table S.2: The 10 types of post-expansion sums for  $(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)})$ .

| Type | Count | Name  | Formula  | Abs.      | Mean |
|------|-------|-------|--|-----------|------|
| Ia   | 4     | $Y_1$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$  | 0         |      |
| Ib   | 8     | $Y_2$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$   | 0         |      |
|      | 4     | $Y_3$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}$   | 0         |      |
| Ic   | 8     | $Y_4$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}$                      | $O(np^2)$ |      |
|      | 4     | $Y_5$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1}$                      | 0         |      |
| Id   | 4     | $Y_6$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1}$ | 0         |      |
| IIa  | 4     | $Z_1$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$  | $O(p^2)$  |      |
|      | 2     | $Z_2$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$  | $O(p^2)$  |      |
| IIb  | 8     | $Z_3$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}$   | 0         |      |
|      | 4     | $Z_4$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$   | 0         |      |
| IIc  | 4     | $Z_5$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1}$                      | $O(np^2)$ |      |
|      | 2     | $Z_6$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{S}(\tilde{\mathbf{P}})_{i_4 i_1}$                       | 0         |      |
| IIIa | 4     | $T_1$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1}$  | $O(p^2)$  |      |
| IIIb | 4     | $T_2$ | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1}$   | 0         |      |
| IV   | 1     | $F$   | $\sum_{\substack{i_1, i_2, i_3, i_4 \\ (\text{dist})}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{D}_{i_4 i_1}$  | $O(p^2)$  |      |

Observe that any sums arising from expansions in  $\mathbb{E}[(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)})^2]$  are bounded by  $(2C_P)^8$  multiplied by a convex combination of terms of the form  $|\mathbb{E}[\prod_{\substack{1 \leq k \leq n \\ 1 \leq l \leq p}} E_{kl}^{t_{kl}}]|$  with  $\sum_{\substack{1 \leq k \leq n \\ 1 \leq l \leq p}} t_{kl} \leq 8$ .

This indicates that such sums are uniformly bounded by a constant, which we denote by  $C_b$ .

Consequently, in our examination of  $\mathbb{E}[Y_i^2], \mathbb{E}[Z_j^2], \mathbb{E}[T_k^2], \mathbb{E}[F^2]$ , we can confine our analysis to the expanded summations in which the indices are distinct. This point will be elaborated further in the proof for Type Ia.

We begin by simplifying the structural assumptions in the matrices. Due to symmetry, the entries  $\mathbf{D}_{ij}, \mathbf{W}_{ij}, \mathcal{S}(\tilde{\mathbf{P}})_{ij}$  are non-zero only when either  $i \leq n$  and  $j > n$ , or vice versa. Therefore, for analyzing terms of the form  $\mathbb{E}[a_{i_1 i_2} b_{i_2 i_3} c_{i_3 i_4} d_{i_4 i_1}]$ , we need only consider two cases: (1)  $i_1, i_3 \leq N$  and  $i_2, i_4 > N$ ; and (2)  $i_2, i_4 \leq N$  and  $i_1, i_3 > N$ . Similarly, when evaluating terms such as  $\mathbb{E}[a_{i_1 i_2} b_{i_2 i_3} c_{i_3 i_4} d_{i_4 i_1} a'_{i'_1 i'_2} b'_{i'_2 i'_3} c'_{i'_3 i'_4} d'_{i'_4 i'_1}]$ , we distinguish four cases based on index configurations: (1)  $i_1, i_3, i'_1, i'_3 \leq N$  and  $i_2, i_4, i'_2, i'_4 > N$ ; (2)  $i_1, i_3, i'_2, i'_4 \leq N$  and  $i_2, i_4, i'_1, i'_3 > N$ ; (3)  $i_2, i_4, i'_1, i'_3 \leq N$  and  $i_1, i_3, i'_2, i'_4 > N$ ; and (4)  $i_2, i_4, i'_2, i'_4 \leq N$  and  $i_1, i_3, i'_1, i'_3 > N$ . Furthermore, it is essential to note that the variables  $\mathbf{E}_{ij}$  are mutually independent, a property that leads to the vanishing of many terms in our expansions. This independence is critical in subsequent analysis.

### S.7.1 Type Ia

$$|\mathbb{E}[Y_1]| \leq |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}]| + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}]|$$

$$\begin{aligned}
&= |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{ki_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}]| + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \\
&\quad \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}]| \\
&= 0
\end{aligned}$$

The last equality holds because, for the case  $i_1, i_3 \leq N, i_2, i_4 > N$ , the terms  $\frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{ki_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \mathbf{W}_{i_2 i_3}$ ,  $\mathbf{W}_{i_3 i_4}$ ,  $\mathbf{W}_{i_4 i_1}$  are independent. Similarly, for the case  $i_2, i_4 \leq N, i_3, i_1 > N$ , the terms  $\frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathbf{W}_{i_2 i_3}$ ,  $\mathbf{W}_{i_3 i_4}$  are independent.

Now we turn to  $\mathbb{E}[Y_1^2]$ . To proceed, we partition the post-expansion summations in  $\mathbb{E}[Y_1^2]$  according to the indices involved.

$$\begin{aligned}
\mathbb{E}[Y_1^2] &= \mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ (i'_1, i'_2, i'_3, i'_4) \text{ dist}}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}] \\
&\leq |\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&\quad + |\sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ (i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ \#\{(i_1, i_2, i_3, i_4) \cap (i'_1, i'_2, i'_3, i'_4)\} \geq 1}} \mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]|
\end{aligned}$$

The second term can be controlled by  $CN^7 \cdot C_b = o(N^8)$ , allowing us to focus solely on the first term.

$$\begin{aligned}
&|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&\leq |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{ki_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{ki'_2}}{|C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}|} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&\quad + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_1, i_3, i'_2, i'_4 \leq N, i_2, i_4, i'_1, i'_3 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{ki_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}} \mathbf{W}_{i'_1 l}}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&\quad + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_1, i'_3 \leq N, i_1, i_3, i'_2, i'_4 > N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{ki'_2}}{|C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}|} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&\quad + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}} \mathbf{W}_{i'_1 l}}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]|
\end{aligned}$$

=0

The last equality holds because, for the cases  $i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N$  and  $i_1, i_3, i'_2, i'_4 \leq N, i_2, i_4, i'_1, i'_3 > N$ , the terms  $\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}$ ,  $\mathbf{W}_{i_3 i_4}$ ,  $\mathbf{W}_{i_4 i_1}$  are independent. Similarly, for the case  $i_2, i_4, i'_1, i'_3 \leq N, i_1, i_3, i'_2, i'_4 > N$  and  $i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N$ , the terms  $\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}$ ,  $\mathbf{W}_{i_2 i_3}$ ,  $\mathbf{W}_{i_3 i_4}$  are independent.

Thus we have proved  $\mathbb{E}[Y_1^2] = o(N^8)$ .

### S.7.2 Type Ib

Similar to the analysis for Type Ia above, we obtain

$$|\mathbb{E}[Y_2]| = |\mathbb{E}[Y_3]| = 0,$$

$$\begin{aligned} \mathbb{E}\left[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\tilde{\mathbf{P}})_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] &= 0, \\ \mathbb{E}\left[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] &= 0. \end{aligned}$$

### S.7.3 Type Ic

Using the property that  $\mathbf{E}_{ij}$  are mutually independent again and  $|\mathcal{S}(\tilde{\mathbf{P}})_{ij}| \leq 2C_P$ , we readily obtain

$$\begin{aligned} |\mathbb{E}[Y_4]| &\leq |\mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}\right]| + |\mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}\right]| \\ &= |\mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}\right]| \\ &\quad + |\mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}\right]| \\ &= |\mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1}\right]| \\ &= \left| \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4}\right] \right| \\ &\leq \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbb{E}\left[\left|\frac{\mathbf{W}_{i_1 i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4}\right|\right] \end{aligned}$$

Now we can select  $(i_2, i_4)$  as follows: we first select a pseudo cluster based on  $\widehat{\mathbf{Z}}^{(m,0)}$  and select a pair  $(i_2, i_4)$  from this cluster. By combining this with the moment inequality for sub-exponential distributions,

$$|\mathbb{E}[Y_4]| \leq 4C_P^2 p^2 \sum_{\substack{(i_2, i_4) \text{ dist} \\ i_2, i_4 \leq N \\ i_4 \in C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_2}^2}{|C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}|}\right] \leq 4C_P^2 p^2 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^2}{|C_i^{(m)}|} = O(np^2)$$

Now we turn to finding an upper bound for  $|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]|$ . Similar to the analysis for Type Ic above, we obtain

$$\begin{aligned} & |\mathbb{E}\left[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right]| \\ & \leq \mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N}} \frac{\sum_{k \in C_{\widehat{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}}{|C_{\widehat{\mathbf{z}}_{i_1}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{k \in C_{\widehat{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{k i'_2}}{|C_{\widehat{\mathbf{z}}_{i'_1}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] \\ & + \mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_1, i_3, i'_2, i'_4 \leq N, i_2, i_4, i'_1, i'_3 > N}} \frac{\sum_{k \in C_{\widehat{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}}{|C_{\widehat{\mathbf{z}}_{i_1}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{l \in C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}} \mathbf{W}_{i'_1 l}}{|C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] \\ & + \mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_1, i'_3 \leq N, i_1, i_3, i'_2, i'_4 > N}} \frac{\sum_{l \in C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{k \in C_{\widehat{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{k i'_2}}{|C_{\widehat{\mathbf{z}}_{i'_1}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] \\ & + \mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N}} \frac{\sum_{l \in C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{l \in C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}} \mathbf{W}_{i'_1 l}}{|C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] \\ & = \mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N \\ i_4 \in C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}, i'_4 \in C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_4}^2}{|C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4}\right] \mathbb{E}\left[\frac{\mathbf{W}_{i'_1 i'_4}^2}{|C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4}\right]\right] \\ & \leq \mathbb{E}\left[\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N \\ i_4 \in C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}, i'_4 \in C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[\left|\frac{\mathbf{W}_{i_1 i_4}^2}{|C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i_3 i_4}\right|\right] \mathbb{E}\left[\left|\frac{\mathbf{W}_{i'_1 i'_4}^2}{|C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}|} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\widetilde{\mathbf{P}})_{i'_3 i'_4}\right|\right]\right] \\ & \leq 16C_P^4 \sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N \\ i_4 \in C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}, i'_4 \in C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_4}^2}{|C_{\widehat{\mathbf{z}}_{i_2}}^{(m)}|}\right] \mathbb{E}\left[\frac{\mathbf{W}_{i'_1 i'_4}^2}{|C_{\widehat{\mathbf{z}}_{i'_2}}^{(m)}|}\right] \end{aligned}$$

Now we can add some terms to make this summation more organized,

$$\begin{aligned}
& \mathbb{E} \left[ \sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\tilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1} \right] \\
& \leq 16 C_P^4 \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}}} \mathbb{E} \left[ \frac{\mathbf{W}_{i_1 i_4}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \right] \sum_{\substack{(i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i'_1, i'_3 > N, i'_2, i'_4 \leq N \\ i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E} \left[ \frac{\mathbf{W}_{i'_1 i'_4}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|} \right] \\
& \leq 16 C_P^4 p^4 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^2}{|C_i^{(m)}|} \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^2}{|C_i^{(m)}|} \\
& = O(n^2 p^4)
\end{aligned}$$

On the other hand, it is easy to obtain  $\mathbb{E}[Y_5] = 0$  and  $\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\tilde{\mathbf{P}})_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1}] = 0$ , following a similar analysis as for Type Ia.

#### S.7.4 Type Id

Similar to the analysis for Type Ia above, we obtain  $\mathbb{E}[Y_6] = 0$  and  $\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\tilde{\mathbf{P}})_{i'_2 i'_3} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1}] = 0$ .

#### S.7.5 Type IIa

Using the previously demonstrated proof approach, we can obtain

$$\begin{aligned}
|\mathbb{E}[Z_1]| & \leq \mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \right] + \mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \right] \\
& = |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} W_{k i_2} \sum_{l \in C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}} \mathbf{W}_{i_2 l}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \frac{\mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}}{|C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}|} \right] \\
& \quad + |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l} \sum_{k \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{k i_3}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \frac{\mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \right] \\
& = |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l} \sum_{k \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{k i_3}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \frac{\mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \right]|
\end{aligned}$$

Again, we can limit our focus on the case where  $i_2$  and  $i_4$  belong to the same pseudo cluster. Hence we obtain

$$|\mathbb{E}[Z_1]| \leq p^2 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} = O(\sigma^4 p^2)$$

Now we turn to finding an upper bound for  $|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]|$ .

$$\begin{aligned}
& |\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&= |\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}, i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}[\frac{\mathbf{W}_{i_1 i_4}^2 \mathbf{W}_{i_3 i_4}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|^2}] \mathbb{E}[\frac{\mathbf{W}_{i'_1 i'_4}^2 \mathbf{W}_{i'_3 i'_4}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|^2}]| \\
&\leq \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}}} \mathbb{E}[\frac{\mathbf{W}_{i_1 i_4}^2 \mathbf{W}_{i_3 i_4}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|^2}] \sum_{\substack{(i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i'_1, i'_3 > N, i'_2, i'_4 \leq N \\ i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}[\frac{\mathbf{W}_{i'_1 i'_4}^2 \mathbf{W}_{i'_3 i'_4}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|^2}] \\
&\leq p^4 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \\
&= O(p^4)
\end{aligned}$$

Similarly, we have

$$\begin{aligned}
|\mathbb{E}[Z_2]| &\leq |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1}]| + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1}]| \\
&= |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \mathbf{W}_{i_2 i_3} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}} \mathbf{W}_{k i_4}}{|C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}|} \mathbf{W}_{i_4 i_1}]| \\
&\quad + |\mathbb{E}[\sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \mathbf{W}_{i_2 i_3} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_4}}^{(m)}} \mathbf{W}_{i_3 l}}{|C_{\tilde{\mathbf{z}}_{i_4}}^{(m)}|} \mathbf{W}_{i_4 i_1}]| \\
&\leq 2p^2 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \\
&= O(p^2)
\end{aligned}$$

Now we turn to finding an upper bound for  $|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]|$ , following a similar analysis as for Type Ia.

$$\begin{aligned}
& |\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\
&= |\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N \\ i_3 \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}, i'_3 \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}}} \mathbb{E}[\frac{\mathbf{W}_{i_1 i_4}^2 \mathbf{W}_{i_3 i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|^2}] \mathbb{E}[\frac{\mathbf{W}_{i'_1 i'_4}^2 \mathbf{W}_{i'_3 i'_2}^2}{|C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}|^2}]|
\end{aligned}$$

$$\begin{aligned}
& + \left| \sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_1, i_3, i'_2, i'_4 \leq N, i_2, i_4, i'_1, i'_3 > N \\ i_3 \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}, i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[ \frac{\mathbf{W}_{i_1 i_4}^2 \mathbf{W}_{i_3 i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|^2} \right] \mathbb{E}\left[ \frac{\mathbf{W}_{i'_1 i'_4}^2 \mathbf{W}_{i'_3 i'_2}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|^2} \right] \right| \\
& + \left| \sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_2, i_4, i'_1, i'_3 \leq N, i_1, i_3, i'_2, i'_4 > N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}, i'_3 \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}}} \mathbb{E}\left[ \frac{\mathbf{W}_{i_1 i_4}^2 \mathbf{W}_{i_3 i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|^2} \right] \mathbb{E}\left[ \frac{W_{i'_1 i'_4}^2 \mathbf{W}_{i'_3 i'_2}^2}{|C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}|^2} \right] \right| \\
& + \left| \sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}, i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[ \frac{\mathbf{W}_{i_1 i_4}^2 \mathbf{W}_{i_3 i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|^2} \right] \mathbb{E}\left[ \frac{\mathbf{W}_{i'_1 i'_4}^2 \mathbf{W}_{i'_3 i'_2}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|^2} \right] \right| \\
& \leq 4p^4 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \\
& = O(p^4)
\end{aligned}$$

### S.7.6 Type IIb

Similar to the analysis for Type Ia above, we obtain

$$\begin{aligned}
& |\mathbb{E}[Z_3]| = |\mathbb{E}[Z_4]| = 0, \\
& \mathbb{E}\left[ \sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1} \right] = 0, \\
& \mathbb{E}\left[ \sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} D_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathcal{S}(\tilde{\mathbf{P}})_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1} \right] = 0.
\end{aligned}$$

### S.7.7 Type IIc

Using the previously demonstrated proof approach, we can obtain

$$\begin{aligned}
|\mathbb{E}[Z_5]| & \leq |\mathbb{E}\left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \right]| + |\mathbb{E}\left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \right]| \\
& = |\mathbb{E}\left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}} \mathbf{W}_{i_2 l}}{|C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \right]| \\
& \quad + |\mathbb{E}\left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l}}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|} \frac{\sum_{k \in C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}} \mathbf{W}_{k i_3}}{|C_{\tilde{\mathbf{z}}_{i_3}}^{(m)}|} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \right]|
\end{aligned}$$

We can still limit our focus on the case where  $i_1$  and  $i_3$  belong to the same pseudo cluster. However, unlike before, this time there are more non-zero terms in the numerator.

$$|\mathbb{E}[Z_5]| \leq 4C_P^2 p^2 \mathbb{E} \left[ \sum_{\substack{(i_1, i_3) \text{ dist} \\ i_1, i_3 \leq N \\ i_3 \in C_{\hat{\mathbf{z}}_{i_1}}^{(m)}}} \frac{|C_{\hat{\mathbf{z}}_{i_1}}^{(m)}| \mathbf{W}_{i_1 i_2}^2}{|C_{\hat{\mathbf{z}}_{i_1}}^{(m)}|^2} \right] \leq 4C_P^2 p^2 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^2}{|C_i^{(m)}|} = O(np^2)$$

Now we turn to finding an upper bound for  $|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1}]|$ .

$$\begin{aligned} & |\mathbb{E} \left[ \sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1} \right]| \\ &= \left| \sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N \\ i_3 \in C_{\hat{\mathbf{z}}_{i_1}}^{(m)}, i'_3 \in C_{\hat{\mathbf{z}}_{i'_1}}^{(m)}}} \mathbb{E} \left[ \frac{(\sum_{k \in C_{\hat{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2})^2}{|C_{\hat{\mathbf{z}}_{i_1}}^{(m)}|^2} \mathcal{S}(\tilde{\mathbf{P}})_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \right] \mathbb{E} \left[ \frac{(\sum_{k \in C_{\hat{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{k i'_2})^2}{|C_{\hat{\mathbf{z}}_{i'_1}}^{(m)}|^2} \mathcal{S}(\tilde{\mathbf{P}})_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1} \right] \right| \\ &\leq 16C_P^4 \sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist} \\ i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N \\ i_3 \in C_{\hat{\mathbf{z}}_{i_1}}^{(m)}, i'_3 \in C_{\hat{\mathbf{z}}_{i'_1}}^{(m)}}} \mathbb{E} \left[ \frac{\sum_{k \in C_{\hat{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}^2}{|C_{\hat{\mathbf{z}}_{i_1}}^{(m)}|^2} \right] \mathbb{E} \left[ \frac{\sum_{k \in C_{\hat{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{k i'_2}^2}{|C_{\hat{\mathbf{z}}_{i'_1}}^{(m)}|^2} \right] \\ &= 16C_P^4 p^4 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^2}{|C_i^{(m)}|} \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^2}{|C_i^{(m)}|} \\ &= O(n^2 p^4) \end{aligned}$$

On the other hand, it is easy to obtain  $|\mathbb{E}[Z_6]| = 0$  and  $\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} D_{i_1 i_2} \mathcal{S}(\tilde{\mathbf{P}})_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{S}(\tilde{\mathbf{P}})_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1}] = 0$ , following a similar analysis as for Type Ia.

### S.7.8 Type IIIa

Using the previously demonstrated proof approach, we can obtain

$$\begin{aligned} |\mathbb{E}[T_1]| &\leq |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \right]| + |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \right]| \\ &= |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 \leq N, i_2, i_4 > N}} \frac{\sum_{k \in C_{\hat{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2} \sum_{l \in C_{\hat{\mathbf{z}}_{i_3}}^{(m)}} W_{i_2 l} \sum_{k \in C_{\hat{\mathbf{z}}_{i_3}}^{(m)}} \mathbf{W}_{k i_4}}{|C_{\hat{\mathbf{z}}_{i_1}}^{(m)}|} \mathbf{W}_{i_4 i_1} \right]| \\ &\quad + |\mathbb{E} \left[ \sum_{\substack{(i_1, i_2, i_3, i_4) \text{ dist} \\ i_1, i_3 > N, i_2, i_4 \leq N}} \frac{\sum_{l \in C_{\hat{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{i_1 l} \sum_{k \in C_{\hat{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{k i_3} \sum_{l \in C_{\hat{\mathbf{z}}_{i_4}}^{(m)}} \mathbf{W}_{i_3 l}}{|C_{\hat{\mathbf{z}}_{i_2}}^{(m)}|} \mathbf{W}_{i_4 i_1} \right]| \end{aligned}$$

$$\begin{aligned} &\leq 2p^2 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \\ &= O(p^2) \end{aligned}$$

Here, we still rely on the key observation that two indices not exceeding  $n$  must belong to the same pseudo-cluster to ensure that the corresponding post-expansion is non-zero.

Now we turn to finding an upper bound for  $|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]|$ .

$$\begin{aligned} &|\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]| \\ &= |\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_1, i_3, i'_1, i'_3 \leq N, i_2, i_4, i'_2, i'_4 > N \\ i_3 \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}, i'_3 \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|^3}\right] \mathbb{E}\left[\frac{\mathbf{W}_{i'_1 i'_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{k i'_2}^2}{|C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}|^3}\right]| \\ &+ |\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_1, i_3, i'_2, i'_4 \leq N, i_2, i_4, i'_1, i'_3 > N \\ i_3 \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}, i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}} \mathbf{W}_{k i_2}^2}{|C_{\tilde{\mathbf{z}}_{i_1}}^{(m)}|^3}\right] \mathbb{E}\left[\frac{\mathbf{W}_{i'_1 i'_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}} \mathbf{W}_{k i'_3}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|^3}\right]| \\ &+ |\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_2, i_4, i'_1, i'_3 \leq N, i_1, i_3, i'_2, i'_4 > N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}, i'_3 \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{k i_3}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|^3}\right] \mathbb{E}\left[\frac{\mathbf{W}_{i'_1 i'_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}} \mathbf{W}_{k i'_2}^2}{|C_{\tilde{\mathbf{z}}_{i'_1}}^{(m)}|^3}\right]| \\ &+ |\sum_{\substack{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist} \\ i_2, i_4, i'_2, i'_4 \leq N, i_1, i_3, i'_1, i'_3 > N \\ i_4 \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}, i'_4 \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}}} \mathbb{E}\left[\frac{\mathbf{W}_{i_1 i_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}} \mathbf{W}_{k i_3}^2}{|C_{\tilde{\mathbf{z}}_{i_2}}^{(m)}|^3}\right] \mathbb{E}\left[\frac{\mathbf{W}_{i'_1 i'_4}^2 \sum_{k \in C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}} \mathbf{W}_{k i'_3}^2}{|C_{\tilde{\mathbf{z}}_{i'_2}}^{(m)}|^3}\right]| \\ &\leq 4p^4 \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \sum_{i=1}^m \binom{|C_i^{(m)}|}{2} \frac{C\sigma^4}{|C_i^{(m)}|^2} \\ &= O(p^4) \end{aligned}$$

### S.7.9 Type IIIb

Similar to the analysis for Type Ia above, we obtain  $|\mathbb{E}[T_2]| = 0$  and  $\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1}] = 0$ .

### S.7.10 Type IV

Similar to the analysis for Type IIIa above, we obtain  $|\mathbb{E}[F]| = O(p^2)$  and  $\mathbb{E}[\sum_{(i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4) \text{ dist}} \mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathcal{S}(\tilde{\mathbf{P}})_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathcal{S}(\tilde{\mathbf{P}})_{i'_4 i'_1}] = O(p^4)$ .

Thus, we have proved  $|\mathbb{E}[Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}]| = O(\sigma^4 p^2)$ ,  $\text{Var}(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)}) = o(N^8)$ .

Note that if  $m = K$ , then  $\mathcal{S}(\tilde{\mathbf{P}})$  reduces to a zero matrix. Thus, any post-expansion summand that involves  $\mathcal{S}(\tilde{\mathbf{P}})$  is zero. Then it follows that

$$Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)} = 4Y_1 + 4Z_1 + 2Z_2 + 4T_1 + F,$$

indicating that it suffices to analyze  $\mathbb{E}[Y_1^2], \mathbb{E}[Z_1^2], \mathbb{E}[Z_2^2], \mathbb{E}[T_1^2], \mathbb{E}[F^2]$ . Based on preceding results, we now focus on terms where the indices are not distinct.

Each post-expansion term is a convex combination of expressions of the form  $|\mathbb{E}[\prod_{\substack{1 \leq k \leq n \\ 1 \leq l \leq p}} \mathbf{E}_{kl}^{t_{kl}}]|$  with  $\sum_{\substack{1 \leq k \leq n \\ 1 \leq l \leq p}} t_{kl} = 8$ . A nonzero contribution occurs only if  $1 \notin \{t_{kl}\}_{\substack{1 \leq k \leq n \\ 1 \leq l \leq p}}$ , which implies that each term contains at most 4 distinct entries of the matrix  $\mathbf{E}$ . Therefore, it can be seen that the contributions of most of the terms are zero. This is a key observation. The main objective moving forward is to provide the upper bound on the number of non-zero contributing terms in expanded forms of  $\mathbb{E}[Y_1^2], \mathbb{E}[Z_1^2], \mathbb{E}[Z_2^2], \mathbb{E}[T_1^2], \mathbb{E}[F^2]$  respectively.

As a result, to ensure

$$\begin{aligned} & \mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}] \\ &= \mathbb{E}\left[\frac{\sum_{k \in C_{\tilde{z}_{i_1}}^{(m)}} \mathbf{W}_{ki_2}}{|C_{\tilde{z}_{i_1}}^{(m)}|} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \frac{\sum_{k \in C_{\tilde{z}_{i'_1}}^{(m)}} \mathbf{W}_{ki'_2}}{|C_{\tilde{z}_{i'_1}}^{(m)}|} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}\right] \\ &\neq 0, \end{aligned}$$

it is necessary that  $\#\{\mathbf{W}_{i_2 i_3}, \mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}, \mathbf{W}_{i'_2 i'_3}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}\} \leq 4$ . Since  $\mathbf{W}_{i_2 i_3}, \mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}$  and  $\mathbf{W}_{i'_2 i'_3}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}$  are distinct within their respective groups, we must also have

$$\#\{i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4\} \leq 5.$$

Furthermore, although numerous terms are present in  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i'_1 i'_2}$ , only a subset of these terms can contribute to  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]$ . This constraint arises because the power of any  $\mathbf{E}_{ij}$  in the expression must exceed 2. Therefore, for terms in the post-expansion of  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i'_1 i'_2}$  to contribute, they must either satisfy this condition inherently or appear within the set  $\{\mathbf{W}_{i_2 i_3}, \mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}, \mathbf{W}_{i'_2 i'_3}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}\}$ . Consequently, if  $i_1, i'_1 \leq n$ , for example, at most  $|C_{\tilde{z}_{i'_1}}^{(m)}| \mathbb{I}(z_{i_1} = z_{i'_1}) + \binom{4}{2}$  terms can contribute. The result holds analogously in other cases.

By Assumption 2 and Theorem 1, we have  $|C_i^{(m)}| \geq \alpha_0 n$  for  $i = 1, \dots, m$ . Therefore, it follows that  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}] = O(\frac{1}{N})$ . Consequently, we obtain  $\mathbb{E}[Y_1^2] = O(N^4)$ .

Next, we analyze  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]$ .

If  $\#\{\mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}\} = 4$ , then at most one term in the post-expansion of  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3}$  can contribute. Thus,  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{W}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{W}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}] = O(\frac{1}{N^4})$ , and we obtain  $\mathbb{E}[Z_1^2] = N^7 \cdot O(\frac{1}{N^4}) + O(N^2) = O(N^3)$ .

If  $\#\{\mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}\} = 3$ , then exactly two of these terms have a power of 1. In this case, on the order of  $N$  terms in the post-expansion of  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3}$  can contribute,

leading to  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}] = O(\frac{1}{N^3})$ . Consequently,  $\mathbb{E}[Z_1^2] = N^7 \cdot O(\frac{1}{N^3}) + O(N^2) = O(N^4)$ .

If  $\#\{\mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}\} = 2$ , then  $\#\{i_1, i_2, i_3, i_4, i'_1, i'_2, i'_3, i'_4\} \leq 5$ . Here, terms in the expansion of  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3}$  need only satisfy the condition that powers of  $E_{ij}$  in them are not 1, without requiring to ensure powers of  $\{\mathbf{W}_{i_3 i_4}, \mathbf{W}_{i_4 i_1}, \mathbf{W}_{i'_3 i'_4}, \mathbf{W}_{i'_4 i'_1}\}$  are not 1. Therefore, on the order of  $N^2$  terms in the post-expansion of  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3}$  can contribute, yielding  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{W}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{W}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}] = O(\frac{1}{N^2})$ . Thus,  $\mathbb{E}[Z_1^2] = N^5 \cdot O(\frac{1}{N^2}) + O(N^2) = O(N^3)$ .

Combining these results, we conclude that  $\mathbb{E}[Z_1^2] = O(N^4)$ .

Note that our proof does not depend on the order of  $\mathbf{D}$  and  $\mathbf{W}$ ; therefore, we also have  $\mathbb{E}[Z_2^2] = O(N^4)$ .

Finally, we analyze  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{W}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{W}_{i'_4 i'_1}]$  and  $\mathbb{E}[\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{D}_{i'_4 i'_1}]$  using a similar approach. Although both  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4}$  and  $\mathbf{D}_{i_1 i_2} \mathbf{D}_{i_2 i_3} \mathbf{D}_{i_3 i_4} \mathbf{D}_{i_4 i_1} \mathbf{D}_{i'_1 i'_2} \mathbf{D}_{i'_2 i'_3} \mathbf{D}_{i'_3 i'_4} \mathbf{D}_{i'_4 i'_1}$  contain numerous terms, it is easy to verify only on the order of at most  $N^3$  and  $N^4$  terms, respectively, can contribute. Since we now only need to consider post-expansion terms in  $\mathbb{E}[T_1^2]$  and  $\mathbb{E}[F^2]$  where indices are not distinct, and there are only on the order of  $N^7$  such terms, it follows that  $\mathbb{E}[F^2] = O(N^4)$ .

By combining the means of these terms, we deduce that  $\mathbb{E}[(Q_N^{(m,0)} - \tilde{Q}_N^{(m,0)})^2] = O(N^4)$ .

This concludes the proof of Lemmas S.2 and S.5.

## S.8 Proof of Theorem 1

We need two main theorems.

**Theorem S.1.** *Consider the settings and assumptions in Section 3, and suppose  $\mathbf{P} = \mathbf{U}\Sigma\mathbf{V}^\top$ ,  $\mathbf{X} = \mathbf{P} + \mathbf{E} = \mathbf{U}_\mathbf{X}\Sigma_\mathbf{X}\mathbf{V}_\mathbf{X}^\top$ . We define  $\mathbf{H}_\mathbf{U} := \mathbf{U}_\mathbf{X}^\top \mathbf{U}$  and  $\mathbf{H}_\mathbf{V} := \mathbf{V}_\mathbf{X}^\top \mathbf{V}$ . With probability at least  $1 - O(N^{-5})$ , one has*

$$\max \left\{ \|\mathbf{U}_\mathbf{X} \text{sgn}(\mathbf{H}_\mathbf{U}) - \mathbf{U}\|_{2,\infty}, \|\mathbf{V}_\mathbf{X} \text{sgn}(\mathbf{H}_\mathbf{V}) - \mathbf{V}\|_{2,\infty} \right\} \lesssim \frac{\kappa(\mathbf{P})\sqrt{N}}{\sigma_K}, \quad (\text{S.2})$$

provided that  $\sigma_K = \omega(\kappa(\mathbf{P})\sigma\sqrt{N})$

*Proof.* Since each  $\mathbf{E}_{ij}$  is now generated from a sub-exponential distribution, it follows directly that there exists an event  $B_N$  such that  $\mathbb{P}(B_N^c) \leq N^2 \exp\left(-\frac{N^{0.4}}{C\sigma}\right)$  as  $N \rightarrow \infty$ , and on the event  $B_N$ , we have  $|\mathbf{E}_{ij}| \leq N^{0.4}$ . Consequently, we can replicate the analysis employed in the proof of Theorem 4.4 in Chen et al. (2021) to establish the desired result.  $\square$

**Remark 9.** *In the proof of Chen et al. (2021), the authors assume that  $|\mathbf{E}_{ij}| \leq B$ , where  $B = O\left(\sigma\sqrt{\frac{K}{\log(N)}}\|\mathbf{U}\|_{2,\infty}^2\right)$ , which does not align with our setting here. However, by setting  $B = N^{0.4}/\log(N)$  and following a similar line of reasoning, we can establish the desired conclusion. Given the tedious nature of the details, which are almost identical to those presented in Chen et al. (2021), we omit them here.*

**Definition 2** (Distance-based metrics defined by bottom up pruning). *Fixing  $K > 1$  and  $1 < m \leq K$ , consider a  $K \times (m-1)$  matrix  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K]^\top$ . First, let  $d_K(\mathbf{U})$  be the minimum pairwise distance of all  $K$  rows. Second, let  $\mathbf{u}_k$  and  $\mathbf{u}_\ell$  ( $k < \ell$ ) be the pair that satisfies  $\|\mathbf{u}_k - \mathbf{u}_\ell\| = d_K(\mathbf{U})$  (if this holds for multiple pairs, pick the first pair in the lexicographical order). Remove*

row  $\ell$  from the matrix  $\mathbf{U}$  and let  $d_{K-1}(\mathbf{U})$  be the minimum pairwise distance for the remaining  $(K - 1)$  rows. Repeat this step and define  $d_{K-2}(\mathbf{U}), d_{K-3}(\mathbf{U}), \dots, d_2(\mathbf{U})$  recursively. Note that  $d_K(\mathbf{U}) \leq d_{K-1}(\mathbf{U}) \leq \dots \leq d_2(\mathbf{U})$ .

**Theorem S.2** (Theorem 4.1 in Jin et al. (2022)). *Fix  $1 < m \leq K$  and let  $n$  be sufficiently large. Consider the non-stochastic vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$  that take only  $K$  values in  $\mathbf{u}_1, \dots, \mathbf{u}_K$ . Write  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_K]^\top$ . Let  $F_k = \{1 \leq i \leq n : \mathbf{x}_i = \mathbf{u}_k\}$ ,  $1 \leq k \leq K$ . Suppose for some constants  $0 < \alpha_0 < 1$  and  $C_0 > 0$ ,  $\min_{1 \leq k \leq K} |F_k| \geq \alpha_0 n$  and  $\max_{1 \leq k \leq K} \|\mathbf{u}_k\| \leq C_0 \cdot d_m(\mathbf{U})$ . We apply the  $k$ -means clustering to a set of  $n$  points  $\widehat{\mathbf{x}}_1, \widehat{\mathbf{x}}_2, \dots, \widehat{\mathbf{x}}_n$  assuming  $\leq m$  clusters, and denote by  $\widehat{S}_1, \widehat{S}_2, \dots, \widehat{S}_m$  the obtained clusters (if the solution is not unique, pick any of them). There exists a constant  $c > 0$ , which only depends on  $(\alpha_0, C_0, m)$ , such that, if  $\max_{1 \leq i \leq n} \|\widehat{\mathbf{x}}_i - \mathbf{x}_i\| \leq c \cdot d_m(\mathbf{U})$ , then  $\#\{1 \leq j \leq m : \widehat{S}_j \cap F_k \neq \emptyset\} = 1$ , for each  $1 \leq k \leq K$ .*

Now let's return to our original question.

For a matrix  $\mathbf{A}$ ,  $\mathbf{A}_{1:m}$  denotes the first  $m$  columns of  $\mathbf{A}$ . By Theorem S.1, we have the following lemma.

**Lemma S.7.** *As  $N \rightarrow \infty$ , with probability  $1 - O(N^{-5})$ , there exists an orthogonal  $K \times K$  matrix  $\mathbf{O}$  such that  $\|r_i((\mathbf{U}\mathbf{x})_{1:m}) - r_i((\mathbf{U}\mathbf{O})_{1:m})\| \leq \|r_i((\mathbf{U}\mathbf{x})_{1:K}) - r_i((\mathbf{U}\mathbf{O})_{1:K})\| \leq C \frac{\kappa(\mathbf{P})\sqrt{N}}{\lambda_K}$  for each  $1 \leq i \leq n$ .*

Note that the matrix  $\mathbf{P}$  has only  $K$  distinct rows, and similarly,  $\mathbf{U}$  also contains  $K$  distinct rows. Consequently, for each  $1 \leq m \leq K$ , the submatrix  $(\mathbf{U}\mathbf{O})_{1:m}$  consists of at most  $K$  distinct rows. Therefore, we can select  $K$  distinct rows from  $(\mathbf{U}\mathbf{O})_{1:m}$  to construct new matrices, denoted as  $(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}$ . Specifically, the construction process involves first selecting  $K$  distinct rows from  $\mathbf{U}$ , multiplying these rows by  $\mathbf{O}$ , and then extracting the first  $m$  columns of the resulting product.

To prove Theorem 1, we apply Lemma S.7 with  $\mathbf{U} = (\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}$ ,  $\mathbf{x}_i = r_i((\mathbf{U}\mathbf{x})_{1:m})$ , and  $\widehat{\mathbf{x}}_i = r_i((\mathbf{U}\mathbf{x})_{1:m})$ , and the main condition we need is  $c_1 \leq d_m((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$  uniformly for all  $\mathbf{O}$ . This is the following lemma.

**Lemma S.8.** *Fix  $1 \leq m \leq K$ . Then there exists a constant  $C > 0$  such that*

$$\min_{\mathbf{O} \in \mathbf{O}^{K \times K}} \{d_m((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})\} \geq C.$$

*Proof.* Below, we fix  $1 < m \leq K$  and a  $K \times K$  orthogonal matrix  $\mathbf{O}$ , and study  $d_m((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$ .

We apply a bottom up pruning procedure to  $(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}$ . First, we find two rows  $r_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$  and  $r_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$  that attain the minimum pairwise distance (if there is a tie, pick the first pair in the lexicographical order) and change the  $l$ -th row to  $r_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$  (suppose  $k < l$ ). Denote the resulting matrix by  $(\mathbf{U}^{(K-1)}(\mathbf{O}))_{1:m}$ . Next, we consider the rows of  $(\mathbf{U}^{(K-1)}(\mathbf{O}))_{1:m}$  and similarly find two rows attaining the minimum pairwise distance and replace one row by the other. Denote the resulting matrix by  $(\mathbf{U}^{(K-2)}(\mathbf{O}))_{1:m}$ .

We repeat these steps to get a sequence of matrices:

$$(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}, (\mathbf{U}^{(K-1)}(\mathbf{O}))_{1:m}, (\mathbf{U}^{(K-2)}(\mathbf{O}))_{1:m}, \dots, (\mathbf{U}^{(2)}(\mathbf{O}))_{1:m}, (\mathbf{U}^{(1)}(\mathbf{O}))_{1:m},$$

where for each  $1 \leq k \leq K$ ,  $(\mathbf{U}^{(k)}(\mathbf{O}))_{1:m}$  has at most  $k$  distinct rows. Comparing it with the Definition 2, we find that  $(\mathbf{U}^{(k-1)}(\mathbf{O}))_{1:m}$  differs from  $(\mathbf{U}^{(k)}(\mathbf{O}))_{1:m}$  in only 1 row, and the difference on this row is a vector whose Euclidean norm is exactly  $d_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$ . As a result,

$$\|(\mathbf{U}^{(k)}(\mathbf{O}))_{1:m} - (\mathbf{U}^{(k-1)}(\mathbf{O}))_{1:m}\| = d_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}), \quad 2 \leq k \leq K.$$

By triangle inequality and the fact that  $d_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}) \leq d_{k-1}((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m})$ , we have

$$\|(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m} - (\mathbf{U}^{(m-1)}(\mathbf{O}))_{1:m}\| \leq \sum_{k=m}^K d_k((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}) \leq (K-m+1) \cdot d_m((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}).$$

To show the claim, it suffices to show that

$$\|(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m} - (\mathbf{U}^{(m-1)}(\mathbf{O}))_{1:m}\| \geq C.$$

Since  $(\mathbf{U}^{(m-1)}(\mathbf{O}))_{1:m}$  has at most  $m-1$  distinct rows, its rank is at most  $m-1$ . Additionally, since  $(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}^\top (\mathbf{U}^{(K)}(\mathbf{O}))_{1:m} = I_m$ , it follows that  $\sigma_m((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}) = 1$

We now combine the results above and apply Weyl's inequality for singular values [Horn and Johnson \(1985\)](#)[Corollary 7.3.5]. It gives

$$1 \leq \sigma_m((\mathbf{U}^{(K)}(\mathbf{O}))_{1:m}) - \sigma_m((\mathbf{U}^{(m-1)}(\mathbf{O}))_{1:m}) \leq \|(\mathbf{U}^{(K)}(\mathbf{O}))_{1:m} - (\mathbf{U}^{(m-1)}(\mathbf{O}))_{1:m}\|.$$

The claim follows immediately.  $\square$