

博士論文

Statistical Analysis for Pattern Extraction from
Relational Data Matrix

(関係データ行列のパターン抽出に関する統計解析)

渡邊 千紘

Abstract

Relational data matrices are everywhere around us, including user-movie rating data and document-word occurrence data. Pattern extraction from such relational data matrices has been extensively studied in the literature to capture and visualize a latent global structure of a given matrix. Specifically, in this paper, we focus on the two well-known problems of relational data analysis: *biclustering* and *matrix reordering*. In the biclustering problem, we assume that a data matrix contains some homogeneous submatrices, say *biclusters*, and estimate the locations of such biclusters. The matrix reordering problem deals with more general structural patterns than biclusters, and its purpose is to find the optimal row/column permutations for a given matrix with which some latent pattern appears (biclusters are special cases here).

An open problem in the biclustering and matrix reordering is that we need to accept various kinds of assumptions in its procedure, such as the number of biclusters or features to be used for row/column orderings. In practice, however, we do not always know the validity of such assumptions in advance. Therefore, we need some evaluation method for the reliability of the model, features, and estimation results of these problems.

In this dissertation, we propose three new approaches for solving this problem in both biclustering and matrix reordering: evaluations of the number of biclusters, the estimated bicluster structure, and the row and column features used for matrix reordering. The first two methods are based on the statistical hypothesis tests, whereas in the last one we maximize the “goodness” of the row/column features in terms of the reconstruction error of the original matrix by a new neural network model.

First, we develop a statistical test on the number of biclusters in a given data matrix A . For a given hypothetical number of biclusters (K_0, H_0) , we test whether matrix A consists of $K_0 \times H_0$ biclusters or more. The proposed test statistic is based on the largest singular value of the standardized data matrix, and its asymptotic distribution based on the null hypothesis is derived by using random matrix theory. We also give a theoretical guarantee for the proposed test statistic under the alternative hypothesis. Based on these results, we propose an asymptotically valid sequential testing on the number of biclusters.

Second, we construct a test on the estimated bicluster structure that have been selected based on a given data matrix A and a specific loss function. Such data-driven selection of a hypothetical bicluster structure is a natural choice when we have no knowledge about

the latent structure of a matrix A in advance. However, to construct a statistically valid test (i.e., the Type I error is controlled by a given significance rate), we need to take the *selective bias* into account. If we derive the p -value of the test statistic based on an invalid assumption that the hypothetical bicluster structure is independent of the data matrix, the test is biased towards optimistic. To avoid this difficulty, we develop a statistical test based on the framework of *selective inference*, where we derive the null distribution of the test statistic under the condition that the hypothetical bicluster structure is selected based on the data matrix.

Finally, we propose a new neural network model for matrix reordering, which automatically extracts row and column features from a given matrix, which are later used for determining the row/column orderings. To evaluate the goodness of the extracted features, we assume a generative model of a data matrix with an autoencoder-like architecture, and evaluate the features based on the reconstruction error of the original matrix. By using the trained neural network model, we can not only determine the row/column orderings of a given matrix A but also visualize a global structural pattern in the matrix A as the output of the model.

Our results provide a clue for examining the validity of the assumptions used in the relational data analysis, and they can be used as first-step analysis tools for acquiring knowledge about the latent structure of a given data matrix.

Acknowledgments

First and foremost, I would like to express my sincere gratitude to my supervisor Taiji Suzuki for all his support during my doctoral course. I learned a lot from the research discussions with him and his sincere attitude toward research. He always gave me accurate advice, thoroughly checked the mathematical precision of the proofs, and discussed the possible directions of further development of the research. I would also like to thank Kenji Yamanishi for his essential comments on my research presentation, as well as on the vision of what the doctoral research should be.

I am grateful to Ryohei Hisano, Tomonari Sei, and Fumiyasu Komaki as the Ph.D. committee members of my doctoral dissertation for their time and fruitful discussion, which has significantly improved my research work.

I am also thankful to all the current and past members of the 6th laboratory of Mathematical Informatics for giving me the insightful feedbacks on my research presentation.

I also appreciate the cooperation and understanding of my colleagues in NTT Communication Science Laboratories, where I work as a researcher during the Ph.D. course.

Finally, I would like to express my gratitude to my family, Yuki Watanabe, Juri Watanabe, and Sumio Watanabe, for all their educational and mental support.

Contents

1	Introduction	16
1.1	Structural pattern mining in relational data	16
1.2	Estimation and evaluation of structural pattern in relational data matrix . .	18
1.3	Contribution of this dissertation	20
1.4	Organization of this dissertation	23
1.5	Notations	24
2	Preliminaries	26
2.1	Relational data	26
2.1.1	Two-mode relational data	26
2.1.2	One-mode relational data	27
2.1.3	Key properties of two-mode random relational data matrices . . .	27
2.2	Biclustering problem	30
2.2.1	Regular-grid biclustering	30
2.2.2	Submatrix detection and localization	31
2.3	Matrix reordering problem	31
3	Statistical test on the number of biclusters in a latent block model	32
3.1	Introduction	32
3.2	Problem settings	34
3.3	Related works	37
3.4	Main results: Test statistic for determining the set of cluster numbers . . .	39
3.5	Experiments	47
3.5.1	Realizable case: Convergence of test statistic T in law to Tracy-Widom distribution	47
3.5.2	Unrealizable case: Asymptotic behavior of test statistic T	48
3.5.3	Accuracy of the proposed goodness-of-fit test	51
3.5.4	Real data analysis: Congressional Voting Records Data Set	56
3.6	Discussions	60
3.7	Chapter conclusion	62

CONTENTS

3.A	Proof of $\left \tilde{S}_{kh} - S_{kh} \right = O_p\left(\frac{1}{m}\right)$	62
3.B	Proof of $\frac{ \lambda_1 - \hat{\lambda}_1 }{b^{\text{TW}}} = O_p\left(m^{-\frac{1}{2l} + \epsilon}\right)$ for all $\epsilon \in \left(0, \frac{2}{7}\right)$ in realizable case . . .	64
3.C	Proof of $\hat{\sigma}^* = O_p(1)$ in unrealizable case	74
3.D	Proof of the asymptotic ICL in the Bernoulli case	76
3.E	Jarque–Bera test for selecting the cluster numbers	78
4	Statistical test on the estimated bicluster structure of a relational data matrix	80
4.1	Introduction	80
4.2	Problem settings	84
4.2.1	Notations and assumptions on data matrix	84
4.2.2	Clustering algorithm based on squared residue minimization . . .	85
4.3	Main results: Statistical test on the solution of squared residue minimization	87
4.3.1	Null distribution of test statistic T	87
4.3.2	Statistical test based on truncated chi distribution	90
4.3.3	Approximated test based on simulated annealing	92
4.4	Experiments	96
4.4.1	Exact test in realizable case: $(K_0, H_0) = (K, H)$	97
4.4.2	Exact test in unrealizable cases: $K_0 < K$ or $H_0 < H$	98
4.4.3	Approximated test in both realizable and unrealizable cases . . .	101
4.4.4	Approximated test in the realizable case, $(K, H) = (3, 3), (4, 4), (5, 5)$	104
4.5	Discussions	109
4.6	Chapter conclusion	114
4.A	Proof of (4.10) that $E^{(g)} - E^{(g')} \neq O$ for $g, g' \in \mathcal{G}_{K_0 H_0}, g \neq g'$	115
4.B	Proof of (4.23) that $\text{rank}(E^{(\hat{g})}) = np - K_0 H_0$	115
4.C	Proof that the number of mutually different patterns of block structures with exactly $K_0 \times H_0$ blocks is lower bounded by $K_0^{n-K_0} H_0^{p-H_0}$	117
4.D	Proof that T_E and $(\mathbf{u}_E, \mathbf{z}_E)$ are mutually independent	118
4.E	Sensitivity analysis with respect to the cooling schedule of simulated annealing	120
4.F	Application of computationally efficient biclustering algorithm for estimating the cluster memberships	120
4.G	Null distribution of test statistic with unknown variance σ_0^2	123
5	Matrix reordering method for capturing flexible structural patterns in a relational data matrix	131
5.1	Introduction	131
5.2	Related works	133
5.3	Main results: Deep two-way matrix reordering	136
5.4	Experiments	138

CONTENTS

5.4.1	Preliminary experiment using synthetic datasets	139
5.4.2	Comparison with existing matrix reordering methods	141
5.4.3	Experiment using the divorce predictors dataset	144
5.4.4	Experiment using the metropolis traffic census dataset	148
5.5	Discussions	150
5.6	Chapter conclusion	153
5.A	Application of DeepTMR to the statistical tests in Chapters 3 and 4	154
5.A.1	Application of DeepTMR to the asymptotic test on the number of biclusters in Chapter 3	155
5.A.2	Application of DeepTMR to the selective test on the estimated bicluster structure in Chapter 4	157
5.B	Correspondence of the attribute indices with meanings in the divorce predictors dataset	157
5.C	Correspondence of the indices with locations in the metropolis traffic census dataset	162
	Conclusion	165

List of Tables

1	Experimental settings of learning rate η , number of epochs T (the total number of iterations is given by $\text{ceil}[Tnp/ \mathcal{I}]$), regularization hyperparameter λ , number of sets of row and column indices in a mini-batch $ \mathcal{I} $, and number of units in ROWENC, COLUMNENC, and DEC networks, m^{ROWENC} , $m^{\text{COLUMNENC}}$, and m^{DEC} , respectively (from input to output).	141
5.A1	Biclustering accuracy of the DeepTMR-based method and the hierarchical clustering (HC) in the settings of Gaussian, Bernoulli, and Poisson distributions.	155

List of Figures

2.1	Two-mode and one-mode data matrices. The color of each entry indicates its value. One-mode matrices can be further decomposed into two subgroups: either directed (center) or undirected (right) networks.	27
2.2	Approximated probability density function of TW_1 distribution [148]. . .	28
2.3	Bicluster structures which we assume in regular-grid biclustering (left) and submatrix detection (right).	30
3.1	The sequential order for testing row and column cluster numbers. For example, let the blue entry $(4, 3)$ be the null cluster numbers (K, H) . Based on this sequentially ordered test, the given cluster numbers (K_0, H_0) are always unrealizable (that is, at least one of $K > K_0$ or $H > H_0$ holds), until it reaches to (K, H)	36
3.2	Difference between matrices P , \bar{P} , and \hat{P} in an unrealizable case.	43
3.3	Q-Q plot of test statistic T against the TW_1 distribution in the setting of Gaussian case	49
3.4	Q-Q plot of test statistic T against the TW_1 distribution in the setting of Bernoulli case	49
3.5	Q-Q plot of test statistic T against the TW_1 distribution in the setting of Poisson case	49
3.6	Ratio of the number of trials where test statistic $T \geq t(\alpha)$, where $t(\alpha)$ is the α upper quantile of the TW_1 distribution. The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions. The horizontal line shows the number of rows n in the observed matrix.	50
3.7	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the test statistic T . The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions. If the test statistic $D\sqrt{r}$ is larger than the significance level α , then the null hypothesis that T follows the TW_1 distribution is rejected, and otherwise, the null hypothesis is accepted.	50

LIST OF FIGURES

3.8	Mean test statistic T in the unrealizable case for 100 trials. The null row and column cluster numbers are 4 and 3, respectively. The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions. The horizontal line shows the number of rows n in the observed matrix.	51
3.9	Mean test statistic T divided by $n^{\frac{5}{3}}$ in the unrealizable case for 100 trials. The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions.	51
3.10	Examples of null block structures of the Gaussian LBM. 40×30 observed matrices are plotted for 10 different settings of B ($t = 1, \dots, 10$). The rows and columns of matrix A were sorted according to the null clusters.	53
3.11	Examples of null block structures of the Bernoulli LBM.	53
3.12	Examples of null block structures of the Poisson LBM.	53
3.13	Accuracy of the proposed goodness-of-fit test under 10 different settings of block-wise mean B . The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions.	54
3.14	Correlation coefficients between the proposed test statistic T or the ICL and the eigenvalues of matrix $\hat{Z}^\top \hat{Z}$. As for the horizontal axes, index $j \in \{1, \dots, p\}$ corresponds to the j th largest eigenvalue.	55
3.15	Accuracy of the proposed test and the model selection based on the ICL under five different settings of block-wise mean B	56
3.16	Ratios of the trials where each set of cluster numbers (K_0, H_0) was selected by the proposed test. The cyan rectangles show the null set of cluster numbers (K, H)	57
3.17	Ratios of the trials where each set of cluster numbers (K_0, H_0) was selected by the model selection based on the ICL. The cyan rectangles show the null set of cluster numbers (K, H)	58
3.18	The observed data matrix of the Congressional Voting Records Data Set [44] and its estimated block structures. The black and white elements, respectively, show “yea” and “nay.”	59
3.19	The p -value of the proposed test and the ICL for each setting of a hypothetical set of cluster numbers (K_0, H_0) . In the left figure, we plotted the p -values only for the tested settings (i.e., until the null hypothesis was accepted). As for the ICL, we plotted a part of results ($K_0 \leq 16, H_0 \leq 16$) for visibility. The cyan rectangles show the selected sets of cluster numbers.	60
3.20	The ratio of the trials in which each set of row and column cluster numbers was selected by the proposed test and the ICL. We plotted the results only for the settings which were selected at least once by the proposed test or the ICL.	61
3.B1	Definition of matrix Q	66

LIST OF FIGURES

3.E1	Accuracy of the Jarque–Bera test under 10 different settings of block-wise mean B (Gaussian LBM).	79
1	Examples of the null and estimated bicluster structures of the observed data matrix with the size of $(n, p) = (9, 9)$. The rows and columns of the observed matrix were sorted according to their clusters, and the blue lines indicate the cluster memberships. The biclustering algorithms of the proposed and existing methods [89] do not necessarily yield identical bicluster structures with the same observed matrix.	83
2	Examples of the observed data matrices with the size of $(n, p) = (9, 9)$, which are generated based on the different block-wise means. The title of each figure shows the range of the block-wise mean vector μ_0 . The blue lines show the null cluster memberships. For visibility, we plotted the matrices whose rows and columns were sorted according to their null clusters.	99
3	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the proposed test (4.31) on the set of cluster memberships \hat{g} with the minimum squared residue.	100
4	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the naive test (4.37).	100
5	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) tests. The null hypothesis that p -value follows the uniform distribution on $[0, 1]$ is rejected if $D\sqrt{r} > \alpha$, where α is a given significance level.	101
6	The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 . For this experiment, we used the setting of $n = p$	101
7	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the proposed (left) and naive (right) statistical tests. If there were no null (i.e., $\hat{g} = g^{(N)}$) or alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, respectively, then the corresponding points of FPR or TPR would not have been plotted.	102
8	AUC score in the realizable case for the proposed and naive statistical tests.	102
9	TPR in the unrealizable case (i.e., $K_0 < K$ or $H_0 < H$) with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the proposed (left) and naive (right) statistical tests.	103
10	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the approximated version of the proposed test.	105

LIST OF FIGURES

11	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the approximated version of the naive test (4.37).	105
12	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) approximated tests.	106
13	The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the approximated clustering algorithm in Section 4.3.3. For the experiment, we used the setting of $n = p$	106
14	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests.	107
15	AUC score in the realizable case for the approximated version of the proposed and naive statistical tests. If there were no null (i.e., $\hat{g} = g^{(N)}$) or alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, respectively, then the corresponding bars would not have been plotted.	107
16	TPR in the unrealizable case (i.e., $K_0 < K$ or $H_0 < H$) with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests.	108
17	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) approximated tests, where $(K, H) = (3, 3)$ (top), $(4, 4)$ (middle), and $(5, 5)$ (bottom).	110
18	The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$), for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the approximated clustering algorithm in Section 4.3.3; $(K, H) = (3, 3)$ (top), $(4, 4)$ (middle), and $(5, 5)$ (bottom). For experiment, we used the setting of $n = p$	111
19	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests, where $(K, H) = (3, 3)$	112
20	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests, where $(K, H) = (4, 4)$	112
21	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests, where $(K, H) = (5, 5)$	112
22	AUC score in the realizable case for the approximated version of the proposed and naive statistical tests, where $(K, H) = (3, 3)$	113
23	AUC score in the realizable case for the approximated version of the proposed and naive statistical tests, where $(K, H) = (4, 4)$	113

LIST OF FIGURES

24	AUC score in the realizable case for the approximated version of the proposed and naive statistical tests, where $(K, H) = (5, 5)$	113
4.A1	A data matrix A whose squared residue σ^2 is zero with block structure g . .	116
4.C1	Unique representation of row cluster indexing where n rows are clustered into exactly K_0 clusters. It must be noted that the set of cluster membership vectors $g^{(1)}$ that can be represented in this form is a subset of all the possible cluster membership vectors.	118
4.E1	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes and cooling schedules r , which was computed by the approximated version of the proposed test.	121
4.E2	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes and cooling schedules r , which was computed by the approximated version of the naive test (4.37).	121
4.E3	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) approximated tests under the different cooling schedule settings r	122
4.E4	The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and cooling schedule r , where \hat{g} is output by the approximated clustering algorithm in Section 4.3.3. For the experiment, we used the setting of $n = p$	122
4.E5	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests under the different cooling schedule settings r	123
4.F1	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the approximated version of the proposed test based on the biclustering algorithm in [140].	124
4.F2	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the approximated version of the naive test (4.37) based on the biclustering algorithm in [140].	124
4.F3	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) approximated tests based on the biclustering algorithm in [140].	126
4.F4	The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the biclustering algorithm in [140]. For the experiment, we used the setting of $n = p$	126
4.F5	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests based on the biclustering algorithm in [140].	127

LIST OF FIGURES

4.F6	Mean computation time for estimating the cluster memberships \hat{g} based on the proposed SA algorithm and fast biclustering algorithm in [140]. The error bars indicate the sample standard deviation of the results for 1000 trials.	127
1	Matrix reordering problem. Given an observed matrix A (left), the proposed DeepTMR reorders the rows and columns of matrix A such that the reordered input matrix (center) shows a meaningful or interpretable structure. The proposed DeepTMR provides us with the denoised mean matrix of the reordered matrix (right) as the output of a trained network, as well as row/column ordering.	137
2	Model architecture of DeepTMR. Given an observed matrix A , DeepTMR is trained to reconstruct each entry A_{ij} from one-dimensional row and column features, which are extracted from the i th row and the j th column of matrix A . After training the network, we reorder the rows and columns of matrix A based on the row and column features extracted in the middle layer.	137
3	Results of the LBM . Top figures: original matrix \bar{A} , observed matrix A obtained by applying random row-column permutation to \bar{A} , reordered input matrix $A^{(\pi)}$, and reordered output matrix $\hat{A}^{(\pi)}$ (left to right). Bottom figures: Encoded row and column features \mathbf{g} and \mathbf{h} and reordered row and column features $\mathbf{g}^{(\pi)}$ and $\mathbf{h}^{(\pi)}$ (left to right).	142
4	Results of the SPM for matrices \bar{A} , A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors \mathbf{g} , \mathbf{h} , $\mathbf{g}^{(\pi)}$, $\mathbf{h}^{(\pi)}$.	142
5	Results of the GBM for matrices \bar{A} , A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors \mathbf{g} , \mathbf{h} , $\mathbf{g}^{(\pi)}$, $\mathbf{h}^{(\pi)}$.	143
6	Examples of matrix \bar{A} for the DGM with different levels of noise standard deviation.	145
7	Examples of observed matrix A for the DGM with different levels of noise standard deviation.	145
8	Examples of reordered input matrix $A^{(\pi)}$ for the DGM with different levels of noise standard deviation (DeepTMR).	145
9	Examples of reordered input matrix $A^{(\pi)}$ for the DGM with different levels of noise standard deviation (SVD-Rank-One).	146
10	Examples of reordered input matrix $A^{(\pi)}$ for the DGM with different levels of noise standard deviation (SVD-Angle).	146
11	Examples of reordered input matrix $A^{(\pi)}$ for the DGM with different levels of noise standard deviation (MDS).	146
12	Examples of reordered output matrix $\hat{A}^{(\pi)}$ for the DGM with different levels of noise standard deviation (DeepTMR).	147

LIST OF FIGURES

13	Matrix reordering error of the DeepTMR, SVD-Rank-One, SVD-Angle, and MDS. The error bars indicate the sample standard deviations of the results for 10 trials.	147
14	Results of the divorce predictors dataset for matrices A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors g , h , $g^{(\pi)}$, $h^{(\pi)}$	149
15	Results of the divorce predictors dataset with SVD-Rank-One, SVD-Angle, and MDS.	150
16	Results of the metropolis traffic census dataset for matrices A , $A^{(\pi)}$, and $\hat{A}^{(\pi)}$. For visibility, we plotted the cyan lines to show the sections between the sets of 10 rows or columns (i.e., $\{R1, \dots, R25\}$ and $\{C1, \dots, C25\}$ for rows and columns, respectively). Because the matrix size is $(n, p) = (249, 249)$, R25 and C25 contain nine rows and nine columns, respectively. The correspondence of the indices with the locations is given in Appendix 5.C.	151
17	Results of the metropolis traffic census dataset for vectors g , h , $g^{(\pi)}$, and $h^{(\pi)}$	152
18	Results of the metropolis traffic census dataset with SVD-Rank-One.	153
19	Results of the metropolis traffic census dataset with SVD-Angle and MDS.	154
5.A1	Q-Q plot of test statistic T computed with the DeepTMR-based biclustering method (left) and hierarchical clustering (right) against the TW_1 distribution in the setting of Gaussian case	156
5.A2	Q-Q plot of test statistic T computed with the DeepTMR-based biclustering method (left) and hierarchical clustering (right) against the TW_1 distribution in the setting of Bernoulli case	156
5.A3	Q-Q plot of test statistic T computed with the DeepTMR-based biclustering method (left) and hierarchical clustering (right) against the TW_1 distribution in the setting of Poisson case	156
5.A4	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the approximated version of the proposed test based on the biclustering algorithm using DeepTMR.	158
5.A5	Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the approximated version of the naive test (4.37) based on the biclustering algorithm using DeepTMR.	158
5.A6	Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) approximated tests based on the biclustering algorithm using DeepTMR.	159

LIST OF FIGURES

5.A7	The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the biclustering algorithm using DeepTMR. For the experiment, we used the setting of $n = p$	159
5.A8	FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the approximated version of the proposed (left) and naive (right) statistical tests based on the biclustering algorithm using DeepTMR.	160
5.A9	AUC score in the realizable case for the approximated version of the proposed and naive statistical tests based on the biclustering algorithm using DeepTMR. If there were no null (i.e., $\hat{g} = g^{(N)}$) or alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, respectively, then the corresponding bars would not have been plotted.	160
1	Summary of the main contributions of this dissertation.	166

Chapter 1

Introduction

1.1 Structural pattern mining in relational data

Relational data are a kind of data with the form of d -dimensional array, where $d \geq 2$. Each dimension of an array corresponds to generally different object, and each entry indicates the relationship among the d objects. Specifically, in this dissertation, we focus on the matrix data (i.e., $d = 2$). Such relational data exist everywhere around us (see Section 2.1 for specific examples), and there have been many studies for analyzing them. For instance, relational data analysis includes:

- **Biclustering:** We assume that a data matrix contains a homogeneous submatrix or *bicluster*, in which all the entries have similar values [5, 63]. Biclustering is a task to find such biclusters from a given matrix. A well-known model for a regular-grid bicluster structure is a latent block model or an LBM [58], where we assume that each entry in a bicluster independently follows an identical distribution. We give a specific formulation of biclustering in Section 2.2.
- **Community detection in network:** Community detection is similar to biclustering in that its purpose is to find a homogeneous submatrix in an adjacency matrix of a given graph. The difference between these two tasks lies in that, in network community detection, both rows and columns represent the same set of nodes and therefore the submatrix consists of the same set of indices in both row and column directions. In other words, community detection is a task to find the set of node indices among which the edges have similar weights. A probabilistic model in this setting, which corresponds to an LBM in the biclustering setting, is a stochastic block model or an SBM [65]. Aside from an SBM, many studies have proposed community detection algorithms based on the *modularity* [116], which quantifies the assortativity (i.e., there are more intra-community edges than inter-community ones) of a network given the community structure [19, 35, 115].

1. Introduction

- **Matrix reordering or seriation:** Biclustering and community detection can also be seen as a problem to find a set of row and column permutations of a given matrix such that the reordered matrix shows a block structure. In matrix reordering, we consider a more general problem to discover such row and column permutations that reveal some pattern in a given matrix without explicit structural assumption (e.g., bicluster structure) [11, 98]. Such generality is both an advantage and a disadvantage of matrix reordering. While it enables us to capture other structural patterns than block structures, the interpretation of the reordered matrix depends more heavily on the analyst than in the case of biclustering. Matrix reordering methods are called “non-destructive data analysis,” since all the information of entry values of an original data matrix is preserved after applying them [6, 110]. We give a specific formulation of matrix reordering in Section 2.3.
- **Matrix factorization (dimensionality reduction):** Matrix factorization is another approach for knowledge acquisition from a given data matrix $A \in \mathbb{R}^{n \times p}$, where we assume that it can be well approximated by a product of two low-rank matrices $W \in \mathbb{R}^{n \times d}$ and $H \in \mathbb{R}^{d \times p}$: $A \approx WH$ with $d \ll n, p$. This type of relational data analysis includes singular value decomposition [14, 55, 133] and non-negative matrix factorization [85, 86, 87]. By such decomposition, we can interpret each row of the original matrix A as a weighted sum of d templates (i.e., rows of matrix H). Matrix factorization methods can also be used for the purpose of link prediction with a matrix with missing entries [29, 109, 149].
- **Influential node detection based on network centrality:** The main purpose of the above tasks is to capture or reconstruct a global structure of a given relational data matrix. Another direction in network analysis is to detect an important node which would have a significant influence on the other nodes. To date, various different measures for such importance or *centrality* of a node have been proposed (e.g., degree, closeness, and betweenness centralities) [21, 49, 50]. For instance, these centrality measures can be used for identifying a super-spreader of an infectious disease [80, 102, 167] and a key location in the air transportation network [59, 60].

Particularly, in this dissertation, we consider a task to discover a latent structural pattern in a given relational data matrix based on biclustering and matrix reordering. These problems are particularly important in case that we do not have prior knowledge about the relationship between rows or columns of a given matrix. For example, let $A \in \mathbb{R}^{n \times p}$ be a given relational data matrix and each entry A_{ij} represents the rating of the j th item by the i th user. Suppose that we want to develop a recommender system solely based on the data matrix A . A natural assumption would be that similar users like similar items. That is, if the i_1 th and i_2 th users gave high ratings to the same set of items (e.g., mystery novels) and they gave low ratings to the same set of items (e.g., fantasy novels), and if the i_1 th user purchased the j th item and gave it a high rating, we expect that the i_2 th user will also give

1. Introduction

a high rating to the j th item. In practice, this kind of data matrix tends to be sparse, that is, it has many zero entries and few non-zero ones. Therefore, to take the local similarity (i.e., common ratings in small number of items) into account, we need to simultaneously discover the similarity pattern in both users and items. The algorithms of the above two tasks, biclustering and matrix reordering, enable us to find such a homogeneous pattern in a given matrix. By permuting the rows and columns of a data matrix based on the result given by such algorithms, we can visualize and interpret the latent structural pattern in the given matrix.

1.2 Estimation and evaluation of structural pattern in relational data matrix

An open problem in relational data analysis including biclustering and matrix reordering is that we need to accept various kinds of assumptions in its procedure. For instance, in most cases, we assume that a given model (e.g., the number of biclusters) or features are appropriate, which should be fixed in advance. Based on such assumptions, we find an optimal set of row and column permutations of a given matrix. Finally, we check the results and interpret them.

However, in practice, it is not always the case that we know the appropriate model or features to represent the latent structural pattern in a given data matrix in advance. Moreover, an estimation result by a biclustering or matrix reordering method is not always accompanied by a guarantee for its reliability. A research question here is,

How can we evaluate the reliability of the model, features, and estimation results of these problems?

Based on this perspective, there exist roughly two directions of research in these problems: to *estimate* a latent structure of a given matrix based on a predefined model or features and to *evaluate* the given model or features or estimation results. In the subsequent paragraphs, we briefly review the existing studies on the former estimation algorithms and describe open problems with regard to the latter evaluation tasks.

Statistical analysis for biclustering With regard to the biclustering problem, there have been many studies on the estimation algorithms, whereas much less studies have been conducted on the evaluation methods. The estimation algorithms of bicluster structure can be decomposed into several categories, according to their structural constraints (i.e., supposed bicluster arrangements). For instance, we assume that a given matrix consists of a regular-grid or checkerboard bicluster structure (i.e., each row (column) belongs to exactly one row (column) cluster) in some algorithms [33, 34, 82], whereas we consider possibly overlapping bicluster structures in other algorithms [40, 134, 168]. Several studies also assumed a hierarchical bicluster structure [63, 128, 129]. These methods can also be

1. Introduction

categorized based on their computational approaches. For instance, some algorithms seek the (approximately) optimal bicluster assignments in terms of a given loss function (e.g., mean squared residue) [32, 34, 63], while other ones use spectral methods (e.g., singular value decomposition) to avoid directly solving a combinatorial optimization problem [41, 82]. For more comprehensive review, we can refer to [25, 107, 123, 141].

Evaluation in the biclustering problem includes the following problems, on which we focus in this dissertation.

- **Evaluation of the number of biclusters:** As for the number of biclusters in a given matrix, there have been proposed statistical tests and model selection methods based on information criteria and cross-validation (see Section 3.3 for a comprehensive review). The purposes of information criteria (and also cross-validation) and statistical tests are slightly different in that the former methods are used to select the best model from a given set of candidates in terms of some statistical property (e.g., marginal likelihood or generalization error), while the latter ones are used when we would like to explicitly control the probability of Type I error. As we discuss later in Chapter 3, even though there have been many studies on the biclustering algorithm itself, there has been no statistical test on the number of biclusters in a two-mode relational data matrix, where the rows and columns represent generally different objects (we give the definition of two-mode relational data in Section 2.1.1).
- **Evaluation of the estimated bicluster structure:** It is also an important task to evaluate the estimation result of biclustering, as well as its model. For this purpose, several studies have proposed statistical tests on bicluster assignments of rows and columns (see Section 4.1 for a comprehensive review). Since the number of possible bicluster assignments of rows and columns increases in exponential order of matrix size (a proof of this is given in Section 4.C), to test all the patterns of block structure is computationally intractable. Instead, we can select a representative block structure from all the patterns and test it. In case that we do not have any prior knowledge about the block structure of a given matrix, a natural choice is to test the optimal set of block memberships that has been selected by some biclustering algorithm. However, there has been no statistical test with the exact p -value derivation that we can apply in such a problem setting.

Extraction of the row and column features used for matrix reordering In matrix reordering, most studies have proposed a method to reorder the rows and columns of a given matrix based on a predefined features (see Section 5.2 for a comprehensive review). Although these algorithms provide the (approximately) optimal row and column orderings for predefined features, there is no way to evaluate the validity of using such features. As in the case of determining the number of biclusters in the biclustering problem, we do not always have a prior knowledge about what features should be used to capture the structural

1. Introduction

pattern in a given matrix. To solve this problem, it would be useful if we can determine the row and column feature extraction process in a data-driven way.

1.3 Contribution of this dissertation

The main contribution of this dissertation is that we develop three approaches for answering the research question given in Section 1.2. Each approach corresponds to one of the three open problems with regard to the biclustering and matrix reordering tasks described in Section 1.2.

- **Evaluation of the number of biclusters:** We develop a statistical test on the number of biclusters in a given relational data matrix in Chapter 3. By sequentially testing multiple sets of hypothetical row and column cluster numbers (K_0, H_0) in ascending order, we can select the accepted one (\hat{K}, \hat{H}) as the estimated number of biclusters. We derive theoretical guarantee for the proposed test in both *realizable* (i.e., the number of biclusters in a given matrix is equal to (K_0, H_0)) and *unrealizable* (i.e., there are more biclusters than (K_0, H_0) in a given matrix) cases. All the main results here are based on the asymptotic theory in terms of matrix size m , where we consider the limit of $m \rightarrow \infty$. Although there have been several studies for testing the number of communities in a one-mode relational data (i.e., network adjacency matrices) [16, 67, 92], this is the first study that has proposed a test on the number of biclusters for two-mode relational data matrices. To check the effectiveness of the proposed test, we conduct experiments by using both synthetic and practical data matrices. The summary of the characteristics of the proposed method is as follows.
 - We can test and select the number of biclusters in a given matrix.
 - The proposed test is validated in the asymptotic sense in terms of the matrix size. The main theorems are based on the recent results in random matrix theory, which we introduce in Section 2.1.3.
 - We can use an arbitrary biclustering algorithm for estimating the bicluster assignments, as long as it satisfies the consistency condition given in Section 3.2. We can test the number of biclusters in the same procedure regardless of the biclustering algorithm.
- **Evaluation of the estimated bicluster structure:** To evaluate an estimation result of biclustering, in Chapter 4, we construct a statistical test on the estimated bicluster memberships of rows and columns that have been selected based on a given data matrix and a specific loss function. Specifically, the proposed test can be applied to the optimal bicluster assignments in terms of the *squared residue* or the sample variance within the same block [34, 63] (the definition of the squared residue is given

1. Introduction

in Section 4.2.2). Unlike the test in Chapter 3, the validity of this proposed test is shown (i.e., we derived the exact p -value of the proposed test statistic) with a finite size data matrix. A naive approach for constructing such a test would be to assume that the estimated bicluster structure is independent of the data matrix and derive the null distribution of a test statistic. However, this assumption is invalid in our setting, where the estimated bicluster structure is selected based on the data matrix. To overcome this difficulty and construct a valid test on the estimated bicluster structure, we employ a *selective inference* approach [13, 88]. This approach can be used for testing the optimal solution of a problem that can be formulated as a set of quadratic inequalities in terms of the vectorized data matrix, each of whose entry is assumed to follow a Gaussian distribution [103]. Aside from such an exact selective test, we also develop an approximated test based on simulated annealing to reduce the computational cost. The summary of the characteristics of the proposed method is as follows.

- We can test the optimal bicluster structure of a given matrix in terms of the squared residue.
 - The proposed test is validated with a finite size data matrix. The main theorems are based on the theory of selective inference.
 - To reduce the computational complexity, we develop an approximated test based on simulated annealing as well as the exact test.
- **Extraction of the row and column features used for matrix reordering:** To solve the matrix reordering task, we consider a slightly different approach to the biclustering problem. In biclustering, we first assume a specific probabilistic model (i.e., latent block model, which is described later in Section 2.2.1) that depends on the bicluster structure and then seek the optimal bicluster structure that best explains the data matrix. In the spectral and dimension-reduction methods of matrix reordering, which we focus in Chapter 5, we also assume a probabilistic model for a given matrix¹. However, instead of assuming a discrete model that directly depends on the row/column orderings and finding the optimal orderings, we assume a continuous low-dimensional model (e.g., bilinear model in (5.1)) and determine the row/column orderings based on its estimation result. In this case, it is difficult to directly test the row and column orderings. Therefore, we adopt a different approach to “evaluate” the goodness of a matrix reordering result. As the existing spectral and dimension-reduction matrix reordering methods, we develop a generative model of a given data matrix based on a set of row and column features, and evaluate the goodness of the extracted features by the reconstruction error of the original matrix.

¹As for the advantages of these approaches compared to the other ones (e.g., Robinsonian and graph-theoretic methods), we give a detailed explanation in Section 5.2.

1. Introduction

The difference between the proposed and existing methods lies in that the feature extraction, as well as the matrix reconstruction based on the extracted features, is trained as a neural network model by using a data matrix. In other words, the features used for matrix reordering are not predefined, and they are determined in a data-driven way such that they can successfully explain the given data matrix. The summary of the characteristics of the proposed method is as follows.

- The matrix reordering is done by using a new autoencoder-like neural network model, which first extracts the row and column features from a given matrix and then reconstruct each entry of the data matrix based on such extracted features. The encoder (i.e., the shallower part of the neural network which extracts the row and column features) as well as the decoder (i.e., the deeper part of the neural network which reconstructs each entry from the features) are trained by using the data matrix. The extracted row and column features by the trained model can then be used for determining the row and column orderings.
- The proposed model including the feature extraction process is trained via back-propagation to minimize the reconstruction error of the original matrix.
- As we discuss later in Chapter 5, the output of the trained model can be seen as a denoised version of the original data matrix, which provides us with the knowledge of the global structure of the matrix.

The publications contained in or related to this dissertation are as follows.

- (Chapter 3) C. Watanabe and T. Suzuki. “Goodness-of-fit test for latent block models,” *Computational Statistics & Data Analysis*, 154:107090, 2021. doi:10.1016/j.csda.2020.107090.
- (Chapter 4) C. Watanabe and T. Suzuki. “Selective inference for latent block models,” *Electronic Journal of Statistics*, 15(1):3137–3183, 2021. doi:10.1214/21-EJS1853.
- (Chapter 5) C. Watanabe and T. Suzuki. “Deep two-way matrix reordering for relational data analysis,” *Neural Networks*, 146:303–315, 2022. doi:10.1016/j.neunet.2021.11.028.
- (Follow-up work of Chapter 3) C. Watanabe and T. Suzuki. “A goodness-of-fit test on the number of biclusters in a relational data matrix,” arXiv:2102.11658, 2021.
- (Follow-up work of Chapter 5) C. Watanabe and T. Suzuki. “AutoLL: automatic linear layout of graphs based on deep neural network,” 2021 IEEE Symposium Series on Computational Intelligence, 2021.

1.4 Organization of this dissertation

The reminder of this dissertation proceeds as follows.

Chapter 2 First, we describe the formal definitions and formulations of the keywords of this dissertation: one-mode and two-mode relational data (Section 2.1), biclustering (Section 2.2) and matrix reordering (Section 2.3) problems. In Section 2.1.1 and 2.1.2, we also give some examples of such relational data matrices. Moreover, in Section 2.1.3, we introduce some important asymptotic properties of a two-mode random relational data matrix, which we use for deriving our main results in Chapter 3.

Chapter 3 This chapter corresponds to the first contribution of this dissertation (i.e., “**evaluation of the number of biclusters**”) given in Section 1.3, and it is based on the study of [152]. In Section 3.1, we introduce the background of latent block models and their open problem in determining the number of biclusters. Then, in Section 3.2, we formally state the problem and give definitions and assumptions that are necessary for the proofs of the main results. Before stating our main results, we review the related studies and describe the difference between them and the proposed method in Section 3.3. Section 3.4 is the main part of this Chapter, where we derive the asymptotic properties (i.e., Theorems 3.4.1, 3.4.2, and 3.4.3) of the proposed test statistic in both null and alternative cases. These properties give the theoretical guarantee for the proposed sequentially ordered test given in Section 3.2. Next, we demonstrate the effectiveness of the proposed test experimentally in Section 3.5. In this section, we first check the behavior of the proposed test statistic with increasing matrix size and compare it to the theory given in Section 3.4. For reference, we also evaluate the accuracy of the proposed test in terms of model selection and compared it with the existing criterion (although the purpose of the proposed test is not to choose a model from multiple candidates as high accuracy as possible, as we also mention in this section). Aside from the above experiments using synthetic data sets, we apply the proposed test to a practical data and analyze the result. Finally, we discuss the main results of this chapter and refer to their limitations in Section 3.6, and give chapter conclusion in Section 3.7. In the appendices, we show detailed proofs for deriving the results in the previous sections.

Chapter 4 This chapter corresponds to the second contribution of this dissertation (i.e., “**evaluation of the estimated bicluster structure**”) given in Section 1.3, and it is based on the study of [154]. We describe the difficulty in constructing a statistically valid test on the estimated bicluster structure that has been selected based on the data matrix, introduce the framework of selective inference to solve such a problem, and refer to the related studies in Section 4.1. Then, in Section 4.2, we describe the notations and assumptions for the proposed test, including the formulation of the specific biclustering algorithm based on

1. Introduction

squared residue minimization for selecting the bicluster structure of a given matrix. Section 4.3 states the main result of this chapter, that is, the null distribution of the proposed test statistic. We give a proof for the main result and the test procedure (i.e., computation of the p -value). Aside from the exact test, as we mentioned in Section 1.3, we develop an approximated test in Section 4.3.3 to reduce the computational cost. In Section 4.4, we conduct experiments to check the validity of the proposed exact and approximated tests in both realizable and unrealizable cases. Then, we discuss the proposed test and room for improvement in Section 4.5. We conclude this chapter in Section 4.6. In the appendices, we give some proofs for a part of the main result and the additional experimental results.

Chapter 5 This chapter corresponds to the third contribution of this dissertation (i.e., “**extraction of the row and column features used for matrix reordering**”) given in Section 1.3, and it is based on the study of [155]. In Section 5.1, we explain an open problem in matrix reordering and briefly introduce our approach to solve it. Then, in Section 5.2, we review the related studies on matrix reordering problem. Here, we describe three specific algorithms of the spectral and dimension-reduction methods, which are particularly relevant to the proposed approach. The proposed matrix reordering method is described in Section 5.3, including the model formulation and the loss function to be minimized. In Section 5.4, we apply the proposed method to both synthetic and practical data matrices to verify its effectiveness. We analyze the reordering results of several types of data matrices with different structural patterns and compare the accuracy of the proposed method in matrix reordering with that of the existing methods. Finally, in Section 5.5, we discuss the results and possible future directions. The conclusion for this chapter is given in Section 5.6.

Finally, we conclude this dissertation by summarizing it and introducing the follow-up works and future perspectives.

1.5 Notations

Throughout this dissertation, we use the following notations.

Order notations We define the orders of sequence of deterministic variables $\{x_m\}$:

- We denote $x_m = O[f(m)]$ iff there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $Cf(m) \geq |x_m|$ holds.
- We denote $x_m = \Omega[f(m)]$ iff there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $Cf(m) \leq |x_m|$ holds.
- We denote $x_m = \Theta[f(m)]$ iff there exist $C_1, C_2 > 0$ and $M > 0$ such that for all $m \geq M$, $C_1f(m) \leq |x_m| \leq C_2f(m)$ holds.

1. Introduction

Similarly, we also define the probabilistic orders of sequence of random variables $\{X_m\}$:

- We denote $X_m = O_p[f(m)]$ iff for all $\epsilon > 0$, there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $\Pr[Cf(m) \geq |X_m|] \geq 1 - \epsilon$ holds.
- We denote $X_m = \Omega_p[f(m)]$ iff for all $\epsilon > 0$, there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $\Pr[Cf(m) \leq |X_m|] \geq 1 - \epsilon$ holds.
- We denote $X_m = \Theta_p[f(m)]$ iff for all $\epsilon > 0$, there exist $C_1, C_2 > 0$ and $M > 0$ such that for all $m \geq M$, $\Pr[C_1f(m) \leq |X_m| \leq C_2f(m)] \geq 1 - \epsilon$ holds.

Operator and Frobenius norms The operator and Frobenius norms of matrix $A = (A_{ij})_{1 \leq i \leq n, 1 \leq j \leq p} \in \mathbb{R}^{n \times p}$ are given by

$$\|A\|_{\text{op}} = \sup_{\mathbf{u} \in \mathbb{R}^p \setminus \mathbf{0}} \frac{\|A\mathbf{u}\|}{\|\mathbf{u}\|}, \quad \|A\|_{\text{F}} = \sqrt{\sum_{i=1}^n \sum_{j=1}^p A_{ij}^2}, \quad (1.1)$$

respectively.

Chapter 2

Preliminaries

2.1 Relational data

In this dissertation, we use the term *relational data* to indicate matrix-form data $A = (A_{ij})_{1 \leq i \leq n, 1 \leq j \leq p} \in \mathbb{R}^{n \times p}$ with real entries, each of which represents a relationship between two objects. Relational data can be decomposed into two groups: either two-mode or one-mode matrix. In the next subsections, we describe the definition, examples, and several important properties of each of them.

2.1.1 Two-mode relational data

If the rows and columns of a data matrix A represent mutually different objects, we call that A is a two-mode matrix (Figure 2.1 left). In this case, the numbers of rows and columns, n and p , respectively, are not necessarily identical (i.e., A is not necessarily a square matrix). For instance, a matrix A that represents the relationship between tourists and locations is a two-mode relational data. In this example, each i th row indicates the i th tourist and each j th column indicates the j th location. The (i, j) th entry A_{ij} represents the rating of the j th location by the i th tourist. In this dissertation, we mostly consider such two-mode relational data matrices.

The examples of two-mode relational data include the MovieLens datasets that contain movie ratings by users [62], the Jester dataset that contains joke ratings by users [57], the congressional voting dataset [44] that represents relationship between congressmen and their attributes [77, 158], the NeurIPS conference papers dataset that counts words (i.e., columns) in NeurIPS conference papers (i.e., rows) [120], and gene expression datasets that represent expression levels of genes (i.e., rows) under different experimental conditions (i.e., columns) [123, 131].

2. Preliminaries

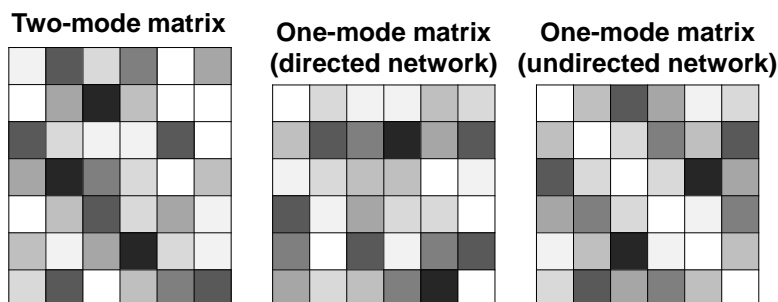


Figure 2.1: Two-mode and one-mode data matrices. The color of each entry indicates its value. One-mode matrices can be further decomposed into two sub-groups: either directed (center) or undirected (right) networks.

2.1.2 One-mode relational data

If the rows and columns of data matrix A represent the same set of objects, we call that A is a one-mode matrix (Figure 2.1 center and right). By definition, the numbers of rows and columns of matrix A are always the same (i.e., $n = p$, A is a square matrix). This kind of data matrix can be viewed as an adjacency matrix of a graph, where the (i, j) th entry A_{ij} indicates the weight of the edge between the i th and the j th nodes. Such adjacency matrices can be further decomposed into two sub-groups: either directed or undirected networks. If the adjacency matrix of a graph is symmetric (i.e., $A_{ij} = A_{ji}$ for all (i, j)), it is called undirected, and otherwise, it is called directed. For instance, an adjacency matrix of a friendship network is a kind of one-mode relational data. Each i th row or column indicates the i th person, and the (i, j) th entry indicates whether the i th and j th persons are friends ($A_{ij} = 1$) or not ($A_{ij} = 0$). In this example, if the network is undirected, we only consider the following two cases: for each pair of persons (i, j) , (1) i and j are friends each other, or (2) i and j are not friends. However, if the network is directed, it corresponds to considering the following four cases: (1) i and j are friends each other, (2) i likes j but j does not like i , (3) j likes i but i does not like j , and (4) i and j do not like each other.

The examples of one-mode relational data (i.e., graph data) include social networks such as Zachary’s karate club network [165] and Facebook network [95], coauthorship networks [93, 113], web graphs [3, 94], a neural network of a creature [156, 157], and a word adjacency network [114].

2.1.3 Key properties of two-mode random relational data matrices

A random matrix is a matrix with random entries. In this dissertation, we assume that a given relational data matrix is a sample of a random matrix, and called it an observed matrix. In this dissertation, we use some recent results with regard to the properties of

2. Preliminaries

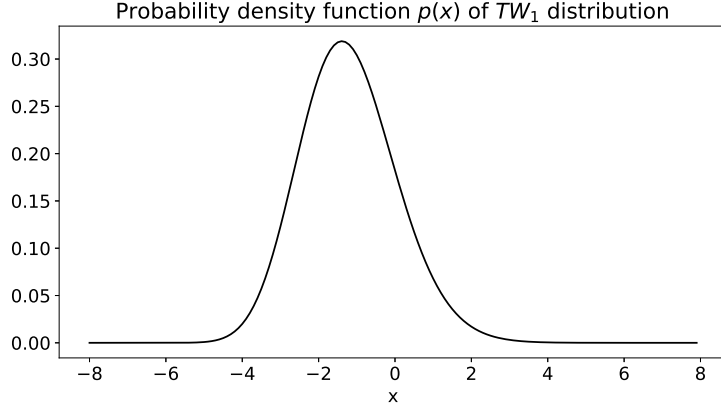


Figure 2.2: Approximated probability density function of TW_1 distribution [148].

two-mode random matrices. Particularly, in Chapter 3, we use the asymptotic distribution of the scaled maximum eigenvalue of a sample covariance matrix and the delocalization property of the eigenvectors of a sample covariance matrix. Before we describe these two properties, we introduce the Tracy-Widom distribution with index 1 (TW_1 distribution), which is used as a null distribution of the proposed test statistic in Chapter 3.

TW_1 distribution The cumulative distribution function of the TW_1 distribution is given by $F_1(x) = E(x)F(x)$, where

$$\begin{aligned} E(x) &= \exp\left(-\frac{1}{2}\int_x^\infty q(y)dy\right), \\ F(x) &= \exp\left(-\frac{1}{2}\int_x^\infty (y-x)q(y)^2dy\right). \end{aligned} \quad (2.1)$$

Here, $q(y)$ is the unique solution of a Painlevé equation of type II (i.e., $\frac{d^2q}{dx^2} = xq + 2q^3$) that satisfies the boundary condition $q(x) \sim \text{Ai}(x)$ in the limit of $x \rightarrow \infty$ with $\text{Ai}(x) = \frac{1}{\pi} \int_0^\infty \cos\left(\frac{t^3}{3} + xt\right) dt$. Since the probability density function of the TW_1 distribution cannot be derived explicitly, an approximated function has been proposed [148], as shown in Figure 2.2.

Asymptotic distribution of the scaled maximum eigenvalue of a sample covariance matrix Let $Z = (Z_{ij})_{1 \leq i \leq n, 1 \leq j \leq p} \in \mathbb{R}^{n \times p}$ be an $n \times p$ matrix, each of whose entries independently follows a (not necessarily identical) distribution with a sub-exponential decay (see Section 3.2 for the precise definition) and with zero mean and unit variance. We assume that the matrix size increases in proportion to some (large) number m (i.e., $n, p \propto m$) and consider the asymptotic property of matrix Z in the limit of $m \rightarrow \infty$. A

2. Preliminaries

recent study [122] revealed that the scaled maximum eigenvalue T^* of sample covariance matrix $Z^\top Z$ converges in law to the TW_1 distribution in the limit of $m \rightarrow \infty$.

Theorem 2.1.1 (Corollary 1.2 of [122]). *Under the above assumptions, in the limit of $m \rightarrow \infty$,*

$$T^* = \frac{\lambda_1 - a^{\text{TW}}}{b^{\text{TW}}} \rightsquigarrow TW_1 \text{ (Convergence in law),} \quad (2.2)$$

where λ_1 is the maximum eigenvalue of the matrix $Z^\top Z$ and

$$a^{\text{TW}} = (\sqrt{n} + \sqrt{p})^2, \quad b^{\text{TW}} = (\sqrt{n} + \sqrt{p}) \left(\frac{1}{\sqrt{n}} + \frac{1}{\sqrt{p}} \right)^{\frac{1}{3}}. \quad (2.3)$$

An important point here is that this convergence in law holds without the assumption that the entries of matrix Z follow an identical distribution. This enables us to utilize this property to prove our main result (Theorem 3.4.1) in Chapter 3, where we define the test statistic based on submatrix-wise standardization of an observed matrix.

Aside from the result in [122], many studies have shown the asymptotic distribution of the above T^* in different problem settings. The studies of [72, 73] have considered a basic setting in which the entries of matrix Z independently follow the standard Gaussian distribution. In [118, 138], the same result has been shown for non-Gaussian cases under some assumptions (e.g., each entry of matrix Z independently follows a symmetric distribution with Gaussian decay and with zero mean and unit variance). The result in [122] is quite general, and it does not require the underlying distribution of each entry of matrix Z to be symmetric. This enables us to apply this result to various types of data matrices, including binary matrices with which we assume that each entry is generated from a Bernoulli distribution.

Delocalization property of the eigenvectors of a sample covariance matrix From Theorem 2.17 of another recent study [17], the eigenvectors of the above sample covariance matrix $Z^\top Z$ satisfy the following *delocalization property*.

Theorem 2.1.2 (Delocalization property of the eigenvectors of matrix $Z^\top Z$ [17]). *Under the above assumptions, the eigenvectors $\{\mathbf{v}_j\}$ ($\|\mathbf{v}_j\| = 1$ for all j) of matrix $Z^\top Z$ satisfy the following property. For all $\tilde{d} \in \mathbb{N}$, for any deterministic vectors $\{\mathbf{w}^{(i)}\}$ that satisfies $\|\mathbf{w}^{(i)}\| = 1$ for $i = 1, \dots, m^{\tilde{d}}$,*

$$\max_{i=1, \dots, m^{\tilde{d}}} \max_{j=1, \dots, p} |\mathbf{v}_j^\top \mathbf{w}^{(i)}| = O_p \left(m^{-\frac{1}{2} + \epsilon} \right), \quad \text{for all } \epsilon > 0. \quad (2.4)$$

2. Preliminaries

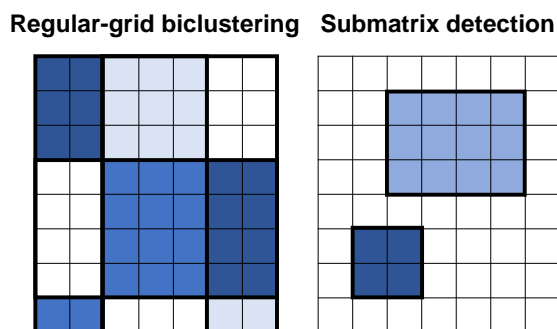


Figure 2.3: Bicluster structures which we assume in regular-grid biclustering (left) and submatrix detection (right).

In Chapter 3, we use this property several times to bound the difference between λ_1 and $\tilde{\lambda}_1$, where $\tilde{\lambda}_1$ is the maximum eigenvalue of matrix $\tilde{Z}^\top \tilde{Z}$ and matrix \tilde{Z} is defined as a random matrix standardized with submatrix-wise **sample** mean and variance (see (3.14) in Section 3.4 for the precise definition). It has also been shown that a similar delocalization property holds for the normalized eigenvectors of a one-mode symmetric random matrix [18].

In this subsection, we introduced two important properties of random matrices without assuming any specific structure behind a data matrix. In the subsequent sections, we formulate two problems, biclustering and matrix reordering, where we assume that a data matrix has some latent structural pattern.

2.2 Biclustering problem

Biclustering is a problem to find a homogeneous submatrix (i.e., bicluster) in a given two-mode relational data matrix A [5, 63]. Specifically, we assume that the entries in a bicluster independently follow an identical distribution. The term of biclustering is used to indicate two different meanings: either clustering all the rows and columns of matrix A to form a regular-grid structure or finding one or more biclusters in matrix A , as shown in Figure 2.3. In Chapters 3 and 4, we consider the former problem setting (i.e., regular-grid structure). In the subsequent subsections, we describe the specific formulations of each setting.

2.2.1 Regular-grid biclustering

Let $A \in \mathbb{R}^{n \times p}$ be an $n \times p$ data matrix. In regular-grid biclustering (Figure 2.3 left), we assume that the set of rows and columns of matrix A are decomposed into row and column clusters, respectively. Specifically, let K and H , respectively, be the numbers of

2. Preliminaries

row and column clusters. We denote the row cluster index of the i th row of matrix A as $g_i^{(1)} \in \{1, \dots, K\}$. Similarly, we denote the column cluster index of the j th column of matrix A as $g_j^{(2)} \in \{1, \dots, H\}$. In regular-grid biclustering, we assume that each entry A_{ij} is independently generated from a bicluster-wise identical distribution. Such a statistical model is called a latent block model (LBM) [58, 63]. In this dissertation, we distinguish an LBM from a stochastic block model (SBM), where we assume that the observed matrix is square symmetric and $g_i^{(1)} = g_i^{(2)}$ holds for all $i \in \{1, \dots, n\}$.

2.2.2 Submatrix detection and localization

Let $A \in \mathbb{R}^{n \times p}$ be an $n \times p$ data matrix. Submatrix detection is a problem to detect the existence of one or more biclusters in matrix A [26, 63, 106, 134]. Particularly, in many studies, the term of submatrix detection is used to indicate the problem to detect large average submatrices, where the mean of the entries is significantly larger than the other entries [26, 27, 101, 106]. Submatrix localization is a problem to estimate the location (i.e., the set of rows and columns) of biclusters in a given matrix A , and it is also called biclustering (Figure 2.3 right). Without the assumption that the entries within a bicluster have larger mean than the other entries, in submatrix localization problem, we can consider more general bicluster structures than in regular-grid biclustering. This is because the $K \times H$ block structure in regular-grid biclustering corresponds to a special case in submatrix localization, where there are $K \times H - 1$ biclusters. Unlike regular-grid biclustering, where each row or column should belong to at least one cluster, we can consider more local bicluster structure in submatrix localization.

2.3 Matrix reordering problem

Matrix reordering is a problem to find a set of row and column permutations $\pi = (\pi^{\text{row}}, \pi^{\text{column}})$ of a given relational data matrix $A \in \mathbb{R}^{n \times p}$ such that the matrix reordered by π shows some structural pattern. Specifically, let π^{row} be a permutation of $\{1, 2, \dots, n\}$ and let π^{column} be a permutation of $\{1, 2, \dots, p\}$. We define that the $\pi^{\text{row}}(i)$ th row ($\pi^{\text{column}}(j)$ th column) in the original matrix corresponds to the i th row (j th column) in the reordered matrix. We denote the reordered matrix as $A^{(\pi)} \in \mathbb{R}^{n \times p}$, each of whose entry is given by $A_{ij}^{(\pi)} = A_{\pi^{\text{row}}(i)\pi^{\text{column}}(j)}$. For a one-mode relational data matrix, the permutations are constrained to satisfy $\pi^{\text{row}} = \pi^{\text{column}}$, and such a task is also called graph reordering or graph layout.

Regular-grid biclustering in Section 2.2.1 can be seen as an example of matrix reordering methods, if we reorder the rows and columns based on their cluster memberships. As we described in Chapter 1, there are various matrix reordering formulations other than the biclustering-based methods, according to the structural pattern that we assume in a given matrix.

Chapter 3

Statistical test on the number of biclusters in a latent block model

Latent block models are used for probabilistic biclustering, which is shown to be an effective method for analyzing various relational data sets. However, there has been no statistical test method for determining the row and column cluster numbers of latent block models. Recent studies have constructed statistical-test-based methods for stochastic block models, which assume that the observed matrix is a square symmetric matrix and that the cluster assignments are the same for rows and columns. In this chapter, we develop a new goodness-of-fit test for latent block models to test whether an observed data matrix fits a given set of row and column cluster numbers, or it consists of more clusters in at least one direction of the row and the column. To construct the test method, we use a result from the random matrix theory for a sample covariance matrix. We experimentally demonstrate the effectiveness of the proposed method by showing the asymptotic behavior of the test statistic and measuring the test accuracy.

3.1 Introduction

Block modeling [5, 63] is known to be effective in representing various relational data sets, such as the data sets of movie ratings [135], customer-product transactions [135], congressional voting [77], document-word relationships [41], and gene expressions [123]. Latent block models or LBMs [58] are used for probabilistic biclustering of such relational data matrices, where rows and columns represent different objects. For instance, suppose that a matrix $A = (A_{ij})_{1 \leq i \leq n, 1 \leq j \leq p} \in \mathbb{R}^{n \times p}$ represents the relationship between users and movies, where entry A_{ij} is the rating of the j th movie by the i th user. In LBMs, we assume a regular-grid block structure behind the observed matrix A ; i.e., both rows (users) and columns (movies) of matrix A are simultaneously decomposed into latent clusters. A block is defined as a combination of row and column clusters, and entries of the same block in

3. Statistical test on the number of biclusters in a latent block model

matrix A are supposed to be i.i.d. random variables.

An open problem in using LBMs is that there has been no statistical procedure for determining the numbers of row and column clusters. Recently, statistical-test-based approaches [16, 67, 92] have been proposed for estimating the cluster number of stochastic block models (SBMs) [65]. SBMs are similar to LBMs in the sense that they assume a block structure behind an observed matrix; however, they are based on different assumptions from LBMs that an observed matrix is a square symmetric matrix and that the cluster assignments are common for rows and columns [108]. In regard to the LBM setting, no statistical method has been constructed to determine row and column cluster numbers. The analysis for an SBM cannot be directly applied to the LBM case due to the difference in the underlying random matrices. To prove the main theorem (i.e., Theorem 3.4.1), we need to derive the orders of multiple variables related to the sample covariance matrix $Z^\top Z$ of a two-mode random matrix Z , such as the number of “large” eigenvalues of $Z^\top Z$ [i.e., t in (3.75)] and an inner product of an eigenvector of $Z^\top Z$ and a deterministic vector. These building blocks of Theorem 3.4.1 should be proven under the different condition from the SBM case.

Aside from the test-based methods, several model selection approaches have been proposed based on cross-validation [30] or an information criterion [76, 77, 119]. However, these approaches have several limitations. (1) First, they cannot provide knowledge about the reliability of the result besides the finally estimated cluster numbers. Rather than minimizing the generalization error, in some cases, it is more appropriate to provide a probabilistic guarantee in reliability for the purpose of knowledge discovery. (2) Second, both the cross-validation-based and information-criterion-based methods depend on the clustering algorithm used. For instance, we can employ the Bayesian information criterion (BIC) for estimating the marginal likelihood only if the Fisher information matrix of the model is regular, which is not the case for block models. Constructing an information criterion that estimates the expectation of the generalization error for a wider class of models is generally difficult. (3) Finally, the above methods require relatively large computational complexity. Computation of an information criterion requires the process of approximating the posterior distribution by the Markov chain Monte Carlo (MCMC) method, and cross-validation requires the iterative calculation of the test error with different sets of partitions of the training and test data sets.

In this chapter, we propose a new statistical test method for LBMs. To construct a hypothesis test with a theoretical guarantee, we use a result from random matrix theory. Recent studies on random matrix theory have revealed the asymptotic behavior of singular values of an $n \times p$ random matrix [7, 8, 43, 54, 72, 73, 118, 122, 137, 138, 161]. Here, we assume that each entry Z_{ij} of matrix Z , which is given by $Z_{ij} = (A_{ij} - P_{ij})/\sigma_{ij}$ (which is computed by the original matrix A , its block-wise mean P and standard deviation σ) follows a distribution with a sub-exponential decay. From the result in [122], the normalized maximum eigenvalue of $Z^\top Z$ converges in law to the Tracy-Widom distribution with index

3. Statistical test on the number of biclusters in a latent block model

1, under the above sub-exponential condition. Based on this result, we construct a goodness-of-fit test for a given set of row and column cluster numbers of an LBM, using the maximum singular value of matrix \hat{Z} , which is an estimator of the matrix Z . We prove that under the null hypothesis (i.e., observed matrix A consists of a given set of row and column cluster numbers), the proposed test statistic T converges in law to the Tracy-Widom distribution with index 1 (Theorem 3.4.1). We also show that under the alternative hypothesis, test statistic T increases in proportion to $m^{\frac{5}{3}}$ with a high probability, where m is a number proportional to the matrix size (Theorems 3.4.2 and 3.4.3).

The proposed method solves the limitations of other model selection approaches. (1) Our statistical test method enables us to obtain knowledge about the reliability of the test results. When testing a given set of row and column cluster numbers, we can explicitly set the probability of Type I error (or false positive) as a significance level α . (2) Unlike the other model selection methods, the proposed method does not depend on the clustering algorithm as long as it satisfies the consistency condition (Section 3.2). It only uses the output of a clustering algorithm to test a given set of cluster numbers; there is no need to modify the test method according to the clustering algorithm. (3) The proposed test method requires relatively small computational complexity. It does not require the MCMC procedure or partitioning into the training and test data sets. For these reasons, the proposed test-based method can be widely used for the purpose of knowledge discovery.

The next sections consist of the detailed explanation of the proposed test method for LBMs. In Section 3.2, we describe the proposed goodness-of-fit test and its theoretical guarantee with the assumptions required for the problem setting. Next, we briefly review the related works and their differences from the proposed method in Section 3.3. The main results are presented in Section 3.4, where we prove the asymptotic properties of the proposed test statistic. In Section 3.5, we experimentally demonstrate the effectiveness of the proposed test method by showing the asymptotic behavior of the test statistic and calculating the test accuracy. We discuss the results and limitations of the proposed method in Section 3.6 and conclude the chapter in Section 3.7.

3.2 Problem settings

Let $A \in \mathbb{R}^{n \times p}$ be an $n \times p$ observed matrix. We assume that each entry of matrix A is independently generated, given its row and column clusters. Let (K, H) be the null set of cluster numbers for rows and columns of an observed matrix A , which is unknown in advance. We denote the cluster indices of the i th row and the j th column of matrix A as $g_i^{(1)} \in \{1, \dots, K\}$ and $g_j^{(2)} \in \{1, \dots, H\}$, respectively. We assume that each entry of matrix A is independently subject to a distribution with the block-wise mean P and the block-wise standard deviation σ :

$$P = (P_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad P_{ij} = B_{g_i^{(1)} g_j^{(2)}}.$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned} \sigma &= (\sigma_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, & \sigma_{ij} &= S_{g_i^{(1)} g_j^{(2)}}. \\ A &= (A_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, & \mathbb{E}[A_{ij}] &= P_{ij}, & \mathbb{E}[(A_{ij} - P_{ij})^2] &= \sigma_{ij}^2, \end{aligned} \quad (3.1)$$

where B_{kh} and $S_{kh} > 0$, respectively, are the mean and the positive standard deviation of entries in the (k, h) th null block under the null hypothesis.

In this chapter, we propose a goodness-of-fit test for selecting the cluster numbers (K, H) from observed matrix A . In such a test, we test whether (K, H) is equal to a given set of cluster numbers (K_0, H_0) or at least one of the given row and column cluster numbers K_0 or H_0 is smaller than the null cluster numbers K or H . In other words, the null (N) and alternative (A) hypotheses are given by

$$(N) : (K, H) = (K_0, H_0), \quad (A) : K > K_0 \text{ or } H > H_0. \quad (3.2)$$

By sequentially testing the cluster numbers in the following order (Figure 3.1), we can select the cluster numbers of a given observed matrix A .

1. Test $(K_0, H_0) = (1, 1)$.
2. Test $(K_0, H_0) = (1, 2), (2, 1)$.
3. Test $(K_0, H_0) = (1, 3), (2, 2), (3, 1)$.
4. ...
5. Test $(K_0, H_0) = (1, L), (2, L-1), \dots, (L, 1)$. Let (\hat{K}, \hat{H}) be the row and column cluster numbers where the null hypothesis is accepted and $\hat{K} + \hat{H} = L + 1$ holds. The selected set of cluster numbers is (\hat{K}, \hat{H}) .

It must be noted that the smaller set of cluster numbers (\hat{K}, \hat{H}) is selected with the smaller significance rate α . Based on the above sequentially ordered test, selection of the cluster numbers requires $(\hat{K} + \hat{H})(\hat{K} + \hat{H} - 1)/2$ tests at most.

Assumptions. Throughout this chapter, we make the following assumptions to derive the test statistics:

- (i). We assume that a distribution of Z_{ij} , which is given by $Z_{ij} = (A_{ij} - P_{ij})/\sigma_{ij}$ as in (3.7) later, has a sub-exponential decay. That is, there exists some $\vartheta > 0$ such that for $x > 1$, $\Pr(|Z_{ij}| > x) \leq \vartheta^{-1} \exp(-x^\vartheta)$. From this assumption, note that for any $\check{n} \in \mathbb{N}$, the \check{n} th moment of a random variable Z_{ij} is finite (i.e., $\mathbb{E}[Z_{ij}^{\check{n}}] < \infty$).
- (ii). We denote the number of rows and columns of matrix A as n and p , respectively. We assume that both n and p increase in proportion to some sufficiently large number m (i.e., $n, p \propto m$).

3. Statistical test on the number of biclusters in a latent block model

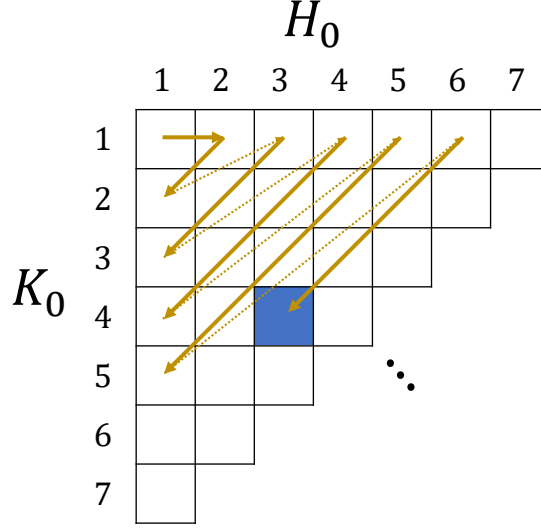


Figure 3.1: The sequential order for testing row and column cluster numbers. For example, let the blue entry $(4, 3)$ be the null cluster numbers (K, H) . Based on this sequentially ordered test, the given cluster numbers (K_0, H_0) are always unrealizable (that is, at least one of $K > K_0$ or $H > H_0$ holds), until it reaches to (K, H) .

- (iii). Let K and H , respectively, be the minimum row and column cluster numbers to represent the block structure of observed matrix A under the null hypothesis. We assume that both K and H are finite constants that do not increase with the matrix sizes n and p . We also assume that the minimum row and column sizes of a block in the null block structure, which we denote as n_{\min} and p_{\min} , respectively, satisfy $n_{\min} = \Omega_p(m)$ and $p_{\min} = \Omega_p(m)$. In other words, we assume that with high probability, there is no “too small” block in matrix A .
- (iv). We assume that the minimum difference between the means of the null blocks satisfies the following conditions.

$$\begin{aligned} \Delta B^{(1)} &= \min_{k, k' \in \{1, \dots, K\}, h \in \{1, \dots, H\}} |B_{kh} - B_{k'h}| = \Omega(1), \\ \Delta B^{(2)} &= \min_{k \in \{1, \dots, K\}, h, h' \in \{1, \dots, H\}} |B_{kh} - B_{kh'}| = \Omega(1). \end{aligned} \quad (3.3)$$

- (v). If the given set of cluster numbers (K_0, H_0) is equal to the null cluster numbers (K, H) , then we call it a *realizable* case. Otherwise, we call it an *unrealizable* case ($K > K_0$ or $H > H_0$). In Section 3.4, we see that Theorems 3.4.2 and 3.4.3 guarantee the behavior of the test statistic T in unrealizable cases. For now, there is

3. Statistical test on the number of biclusters in a latent block model

no way to detect the cases where $(K < K_0) \cap (H \leq H_0)$ or $(K \leq K_0) \cap (H < H_0)$ holds, and to cope with such settings is beyond the scope of this dissertation.

- (vi). In the realizable case, we assume that a clustering algorithm is *consistent*, that is, the probability that it outputs the correct block structure converges to 1, in the limit of $m \rightarrow \infty$. By using this assumption, the proposed method does not depend on a specific clustering algorithm. Several clustering algorithms including [4, 23, 47] have been proven to be consistent. Such consistency in estimating the block structure is also related to the *resolution limit* in network community detection problem for one-mode data matrices [48].

Remark 3.2.1 (Extension of the assumptions with regard to the row and column cluster numbers). *From the result of the follow-up study [153], we can prove the main theorems under more relaxed conditions with regard to the row and column cluster numbers. In the realizable case, Theorem 3.4.1 holds with the following condition:*

$$KH = O\left(m^{\frac{1}{42}-\epsilon_1}\right), \text{ for some } \epsilon_1 > 0. \quad (3.4)$$

$$n_{\min} = \Omega\left(m^{\frac{8}{21}}\right), \quad p_{\min} = \Omega\left(m^{\frac{8}{21}}\right). \quad (3.5)$$

In the unrealizable case, we can prove the upper bound of test statistic T in Theorem 3.4.3 and the lower bound of $T = \Omega_p\left(m^{\frac{2}{3}}\right)$ under the following condition:

$$\frac{KH}{\sqrt{n_{\min}p_{\min}}} = O\left(m^{-\frac{3}{4}-\epsilon_2}\right), \text{ for some } \epsilon_2 > 0. \quad (3.6)$$

3.3 Related works

In this section, we briefly review the related works and explain the differences between them and the proposed method.

Statistical-test-based methods (for SBM) Recently, several methods have been proposed for testing the properties of a given observed matrix in relation to SBMs [16, 67, 74, 92, 164]. Particularly, the methods proposed in [16, 67, 92] have enabled us to estimate the number of blocks for SBMs. However, these methods differ from ours in the problem setting; they can be applied only to an SBM setting, where an observed matrix is a square symmetric matrix, and the cluster assignments are common for rows and columns. There has been no method to estimate the block number for LBMs, where rows and columns (not necessarily square) of an observed matrix are simultaneously decomposed into clusters.

3. Statistical test on the number of biclusters in a latent block model

Cross-validation-based methods Cross-validation is a widely used method for model selection, where a data set is first split into training and test data sets, and then the best model with the minimum test error is determined. Recently, cross-validation methods for matrix data have been proposed [30, 39, 75, 97] to determine the number of clusters in network data. Although the purpose of these methods and our method is similar, these methods differ from ours in that their target is the network data, where the observed matrix is square and its rows and columns represent the same node sets. Thus, the block structure is symmetric regardless of whether the network itself is directed or undirected. Moreover, unlike a statistical test, these methods cannot provide quantitative knowledge about the reliability of the selected model. Furthermore, the computational cost of cross-validation is generally high because it requires the iterative calculation of the test error with different data set partitions.

Information-criterion-based methods Another approach for determining the number of blocks in a matrix is to estimate the generalization error or marginal likelihood by some information criteria for given sets of block numbers. By using such information criteria, we can select a model in a statistically meaningful (non-heuristic) way. In regard to block models, many variants of BIC [66, 76, 77, 119, 130] or MDL [132, 160] have been proposed. Unlike our test-based method, which only requires a clustering algorithm to satisfy the consistency condition (Section 3.2), an information criterion for a theoretical guarantee should be carefully chosen according to the given clustering algorithm. For instance, BIC can be employed for estimating the marginal likelihood only if the Fisher information matrix of the model is regular, which is not the case for block models.

To solve this problem, as an alternative criterion to BIC, the integrated completed likelihood (ICL) criterion has been used in many studies for estimating the number of blocks in LBMs [37, 104, 159]. In ICL, we first derive a marginal likelihood for a given set of an observed matrix and block assignments and then substitute the set of estimated block assignments to approximate the marginal likelihood. However, since ICL is computed based on a single estimator of block assignments, there is no guarantee for the goodness of the approximation of marginal likelihood.

Similar to cross-validation-based methods, information-criterion-based methods cannot provide a probabilistic guarantee for the reliability of the selected model, which is a disadvantage for the purpose of knowledge discovery. The computational cost also becomes a problem because the computation of an information criterion requires the process of approximating the posterior distribution by MCMC.

Other model selection methods Aside from the information criteria, several studies have proposed to determine the number of blocks in LBMs based on the co-clustering adjusted rand index [125], the extended modularity for biclustering [84], or the expected posterior loss for a given loss function [124]. Another approach is to define the posterior

3. Statistical test on the number of biclusters in a latent block model

distribution not only on cluster assignments of rows and columns but also on row and column cluster numbers [117, 158]. Unlike the model selection approaches, such non-parametric Bayesian methods can estimate the distribution of the block numbers. The best-fitted number of the blocks can be determined based on the posterior distribution (e.g., we can choose a MAP estimator [117]). However, in this case, the computational cost of MCMC is higher than that of the information-criterion-based methods because it requires a large number of iterations to approximate the posterior distribution both on the block assignments and the number of blocks.

3.4 Main results: Test statistic for determining the set of cluster numbers

To derive the test statistic for the proposed goodness-of-fit test, we first normalize each entry A_{ij} of an observed matrix A by subtracting P_{ij} and dividing it by σ_{ij} , where P and σ , respectively, are the block-wise mean and standard deviation in (3.1):

$$Z = (Z_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad Z_{ij} = \frac{A_{ij} - P_{ij}}{\sigma_{ij}}. \quad (3.7)$$

By definition, each entry Z_{ij} of matrix Z in (3.7) independently follows a distribution with zero mean and standard deviation of one. Let T^* be the scaled maximum eigenvalue of matrix $Z^\top Z$, which is given by

$$T^* = \frac{\lambda_1 - a^{\text{TW}}}{b^{\text{TW}}}, \quad (3.8)$$

where λ_1 is the maximum eigenvalue of the matrix $Z^\top Z$ and a^{TW} and b^{TW} are defined as in (2.3). According to Theorem 2.1.1 [122], T^* converges in law to the TW_1 distribution in the limit of $m \rightarrow \infty$.

In most cases, the null cluster numbers (K, H) and the null cluster assignments $g^{(1)}$ and $g^{(2)}$ are unknown in advance. Therefore, we can only estimate the block structure based on the observed matrix A and the given cluster numbers. Let (K_0, H_0) be the given set of row and column cluster numbers, and $\hat{g}^{(1)}$ and $\hat{g}^{(2)}$, respectively, be the estimated cluster assignments for rows and columns. Based on such an estimated block structure $(\hat{g}^{(1)}, \hat{g}^{(2)})$, we estimate the block-wise mean and standard deviation by

$$\begin{aligned} \hat{B} &= (\hat{B}_{kh})_{1 \leq k \leq K_0, 1 \leq h \leq H_0}, & \hat{B}_{kh} &= \frac{1}{|I_k||J_h|} \sum_{i \in I_k, j \in J_h} A_{ij}, \\ \hat{P} &= (\hat{P}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, & \hat{P}_{ij} &= \hat{B}_{\hat{g}_i^{(1)} \hat{g}_j^{(2)}}, \\ \hat{S} &= (\hat{S}_{kh})_{1 \leq k \leq K_0, 1 \leq h \leq H_0}, & \hat{S}_{kh} &= \sqrt{\frac{1}{|I_k||J_h|} \sum_{i \in I_k, j \in J_h} (A_{ij} - \hat{P}_{ij})^2}, \end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\hat{\sigma} = (\hat{\sigma}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad \hat{\sigma}_{ij} = \hat{S}_{\hat{g}_i^{(1)} \hat{g}_j^{(2)}}, \quad (3.9)$$

where I_k is the set of row indices of matrix A that are assigned to the k th cluster, and J_h is the set of column indices of matrix A that are assigned to the h th cluster:

$$I_k = \left\{ i : \hat{g}_i^{(1)} = k \right\}, \quad J_h = \left\{ j : \hat{g}_j^{(2)} = h \right\}. \quad (3.10)$$

The consistency assumption (vi) guarantees that if $(K_0, H_0) = (K, H)$, the probability that the cluster assignments $(I_k)_{1 \leq k \leq K_0}$ and $(J_h)_{1 \leq h \leq H_0}$ are correct converges to 1 in the limit of $m \rightarrow \infty$.

We define an estimator of normalized matrix Z in (3.7) based on the estimated block-wise mean \hat{P} and standard deviation $\hat{\sigma}$ in (3.9):

$$\hat{Z} = (\hat{Z}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad \hat{Z}_{ij} = \frac{A_{ij} - \hat{P}_{ij}}{\hat{\sigma}_{ij}}. \quad (3.11)$$

The test statistic T for the proposed goodness-of-fit test is given by the scaled maximum eigenvalue of matrix $\hat{Z}^\top \hat{Z}$:

$$T = \frac{\hat{\lambda}_1 - a^{\text{TW}}}{b^{\text{TW}}}, \quad (3.12)$$

where $\hat{\lambda}_1$ is the maximum eigenvalue of matrix $\hat{Z}^\top \hat{Z}$, and a^{TW} and b^{TW} are given by (2.3).

Based on the following results in Theorems 3.4.1, 3.4.2 and 3.4.3, we propose a one-sided goodness-of-fit test for a given set of cluster numbers (K_0, H_0) at the significance level of α by using the test statistic T :

$$\text{Reject null hypothesis } ((K, H) = (K_0, H_0)), \quad \text{if } T \geq t(\alpha), \quad (3.13)$$

where $t(\alpha)$ is the α upper quantile of the TW_1 distribution. By applying the sequentially ordered test that we explained in Section 3.2 based on the above rejection rule (3.13), we can select a set of row and column cluster numbers (\hat{K}, \hat{H}) for a given observed matrix A .

In the proof of Theorem 3.4.1, we also use the following notations. Let \tilde{B}_{kh} and \tilde{S}_{kh} , respectively, be the **sample** mean and **sample** standard deviation of all the entries in the (k, h) th **null** block in matrix A . Based on such notations, we define the sample mean matrix \tilde{P} and standard deviation matrix $\tilde{\sigma}$ for the correct block structure, and matrix \tilde{Z} by:

$$\begin{aligned} \tilde{P} &= (\tilde{P}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, & \tilde{P}_{ij} &= \tilde{B}_{g_i^{(1)} g_j^{(2)}}, \\ \tilde{\sigma} &= (\tilde{\sigma}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, & \tilde{\sigma}_{ij} &= \tilde{S}_{g_i^{(1)} g_j^{(2)}}, \\ \tilde{Z} &= (\tilde{Z}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, & \tilde{Z}_{ij} &= \frac{A - \tilde{P}_{ij}}{\tilde{\sigma}_{ij}}. \end{aligned} \quad (3.14)$$

3. Statistical test on the number of biclusters in a latent block model

Theorem 3.4.1 (Realizable case). *Under the assumptions in Section 3.2, if $(K_0, H_0) = (K, H)$,*

$$T \rightsquigarrow TW_1 \text{ (Convergence in law),} \quad (3.15)$$

in the limit of $m \rightarrow \infty$, where T is defined as in (3.12).

Proof. First of all, we derive the difference between B_{kh} (S_{kh}) and \tilde{B}_{kh} (\tilde{S}_{kh}), which have been defined in (3.1) and (3.14). Since the number of entries in the block is proportional to m^2 by the assumption (iii), $\sqrt{m^2} (B_{kh} - \tilde{B}_{kh})$ converges to $\mathcal{N}(0, S_{kh}^2)$ from the central limit theorem. Therefore, from Prokhorov's theorem [146], we have

$$\left| \tilde{B}_{kh} - B_{kh} \right| = O_p \left(\frac{1}{m} \right). \quad (3.16)$$

Also, the following equation holds (The proof is given in Appendix 3.A):

$$\left| \tilde{S}_{kh} - S_{kh} \right| = O_p \left(\frac{1}{m} \right). \quad (3.17)$$

From here, we derive the difference between the maximum eigenvalue $\tilde{\lambda}_1$ of the matrix $\tilde{Z}^\top \tilde{Z}$ and the maximum eigenvalue λ_1 of $Z^\top Z$, where the definitions of Z and \tilde{Z} have been given in (3.7) and (3.14), respectively. From Theorem 2.1.1, we have $\lambda_1 = O_p(m)$. Therefore, the largest singular value of Z , which is equal to $\|Z\|_{\text{op}}$, is in the order of $O_p(\sqrt{m})$.

By the subadditivity of the operator norm, we have

$$\left| \|Z\|_{\text{op}} - \|\tilde{Z}\|_{\text{op}} \right| \leq \|Z - \tilde{Z}\|_{\text{op}}. \quad (3.18)$$

Let $A^{(k,h)}$, $P^{(k,h)}$, $\tilde{P}^{(k,h)}$, $Z^{(k,h)}$, and $\tilde{Z}^{(k,h)}$, respectively, be the (k, h) th **null** blocks of matrices A , P , \tilde{P} , Z , and \tilde{Z} . We also denote the row and column sizes of the (k, h) th **null** block as n_k and p_h , respectively. From the definitions in (3.7) and (3.14), we have

$$Z^{(k,h)} = \frac{A^{(k,h)} - P^{(k,h)}}{S_{kh}}, \quad \tilde{Z}^{(k,h)} = \frac{A^{(k,h)} - \tilde{P}^{(k,h)}}{\tilde{S}_{kh}}. \quad (3.19)$$

Combining this with (3.16), (3.17), and the fact that the Frobenius norm upper bounds the operator norm, we have

$$\|Z^{(k,h)} - \tilde{Z}^{(k,h)}\|_{\text{op}} = \left\| \frac{A^{(k,h)} - P^{(k,h)}}{S_{kh}} - \frac{A^{(k,h)} - \tilde{P}^{(k,h)}}{\tilde{S}_{kh}} \right\|_{\text{op}}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&= \left\| \frac{A^{(k,h)} - P^{(k,h)}}{S_{kh}} - \frac{A^{(k,h)} - P^{(k,h)}}{\tilde{S}_{kh}} + \frac{A^{(k,h)} - P^{(k,h)}}{\tilde{S}_{kh}} - \frac{A^{(k,h)} - \tilde{P}^{(k,h)}}{\tilde{S}_{kh}} \right\|_{\text{op}} \\
&\leq \left\| \frac{A^{(k,h)} - P^{(k,h)}}{S_{kh}} - \frac{A^{(k,h)} - P^{(k,h)}}{\tilde{S}_{kh}} \right\|_{\text{op}} + \left\| \frac{A^{(k,h)} - P^{(k,h)}}{\tilde{S}_{kh}} - \frac{A^{(k,h)} - \tilde{P}^{(k,h)}}{\tilde{S}_{kh}} \right\|_{\text{op}} \\
&= \left| \frac{\tilde{S}_{kh} - S_{kh}}{S_{kh}\tilde{S}_{kh}} \right| \|A^{(k,h)} - P^{(k,h)}\|_{\text{op}} + \frac{1}{\tilde{S}_{kh}} \|P^{(k,h)} - \tilde{P}^{(k,h)}\|_{\text{op}} \\
&\leq \left| \frac{\tilde{S}_{kh} - S_{kh}}{S_{kh}\tilde{S}_{kh}} \right| \|A^{(k,h)} - P^{(k,h)}\|_{\text{op}} + \frac{1}{\tilde{S}_{kh}} \|P^{(k,h)} - \tilde{P}^{(k,h)}\|_{\text{F}} \\
&= \left| \frac{\tilde{S}_{kh} - S_{kh}}{\tilde{S}_{kh}} \right| \|Z^{(k,h)}\|_{\text{op}} + \frac{1}{\tilde{S}_{kh}} \sqrt{n_k p_h} |B_{kh} - \tilde{B}_{kh}| \\
&= \frac{O_p(1/m)}{S_{kh} + O_p(1/m)} \|Z^{(k,h)}\|_{\text{op}} + \frac{O_p(1/m)}{S_{kh} + O_p(1/m)} \sqrt{n_k p_h} \quad (\because (3.16), (3.17)) \\
&= \frac{O_p(1/m)}{S_{kh} + O_p(1/m)} O_p(\sqrt{m}) + \frac{O_p(1/m)}{S_{kh} + O_p(1/m)} \sqrt{n_k p_h} \quad (\because \text{Theorem 2.1.1}) \\
&= O_p\left(\frac{1}{\sqrt{m}}\right) + O_p(1) = O_p(1). \tag{3.20}
\end{aligned}$$

Therefore, since the operator norm of a matrix is not larger than the sum of the operator norms of all of its blocks and the number of blocks are finite constants, we have

$$\|Z - \tilde{Z}\|_{\text{op}} \leq \sum_{k=1}^K \sum_{h=1}^H \|Z^{(k,h)} - \tilde{Z}^{(k,h)}\|_{\text{op}} = O_p(1). \tag{3.21}$$

By combining this with (3.18), we obtain

$$\left| \|Z\|_{\text{op}} - \|\tilde{Z}\|_{\text{op}} \right| = O_p(1). \tag{3.22}$$

Next, we consider the joint probability of the event \mathcal{F}_m that $\tilde{Z} = \hat{Z}$ holds and the event $\mathcal{G}_{m,C}$ that $\left| \|Z\|_{\text{op}} - \|\tilde{Z}\|_{\text{op}} \right| \leq C$ holds. Such a joint probability satisfies the following inequality:

$$\Pr(\mathcal{F}_m \cap \mathcal{G}_{m,C}) \geq 1 - \Pr(\mathcal{F}_m^C) - \Pr(\mathcal{G}_{m,C}^C), \tag{3.23}$$

where \mathcal{A}^C is the complement of event \mathcal{A} . The consistency assumption (vi) guarantees that if $(K_0, H_0) = (K, H)$, $\Pr(\mathcal{F}_m^C)$ converges to 0 in the limit of $m \rightarrow \infty$. By combining this fact with (3.22), for all $\epsilon > 0$, there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $\Pr(\mathcal{F}_m \cap \mathcal{G}_{m,C}) \geq 1 - \epsilon$ holds, which results in

$$\left| \|Z\|_{\text{op}} - \|\hat{Z}\|_{\text{op}} \right| = O_p(1). \tag{3.24}$$

3. Statistical test on the number of biclusters in a latent block model

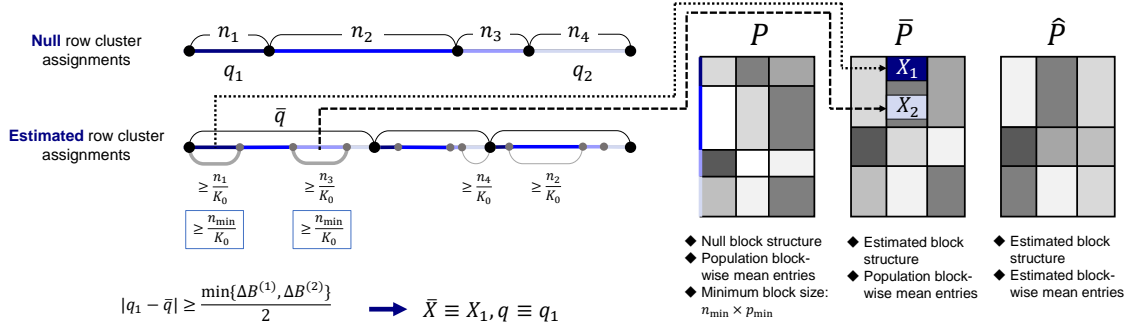


Figure 3.2: Difference between matrices P , \bar{P} , and \hat{P} in an unrealizable case.

By using the above results, we can prove that the following equation holds for all $\epsilon \in (0, \frac{2}{7})$ (The proof is given in Appendix 3.B):

$$\frac{|\lambda_1 - \hat{\lambda}_1|}{b^{\text{TW}}} = O_p \left(m^{-\frac{1}{21} + \epsilon} \right), \quad (3.25)$$

From Theorem 2.1.1, (3.25), and Slutsky's theorem, by setting $\epsilon < \frac{1}{21}$,

$$\frac{\hat{\lambda}_1 - a^{\text{TW}}}{b^{\text{TW}}} = T^* + \frac{\hat{\lambda}_1 - \lambda_1}{b^{\text{TW}}} \rightsquigarrow TW_1 \text{ (Convergence in law)}. \quad (3.26)$$

This is equivalent to the statement of Theorem 3.4.1. \square

Theorem 3.4.2 (Unrealizable case, lower bound). *Under the assumptions in Section 3.2, if $K_0 < K$ or $H_0 < H$,*

$$T = \Omega_p \left(m^{\frac{5}{3}} \right), \quad (3.27)$$

where T is defined as in (3.12).

Proof. Let \bar{P} be a matrix that consists of the **estimated** block structure and whose entries are the population block-wise means, which can be calculated using P (see also Figure 3.2). To derive the difference between matrices P and \hat{P} , we first focus on the relationship between matrices P and \bar{P} . In the unrealizable case (i.e., $K_0 < K$ or $H_0 < H$), we can assume $K_0 < K$ without loss of generality.

Let n_k be the number of rows in the k th **null** row cluster. For all the **null** row cluster indices $k \in \{1, \dots, K\}$, at least one **estimated** row cluster contains n_k/K_0 or more rows that are assigned to the k th row cluster in the **null** block structure (otherwise, the total number of rows in the k th **null** row cluster is smaller than n_k). Since $K_0 < K$, at least one estimated block contains two or more sets of rows whose **null** row clusters are mutually

3. Statistical test on the number of biclusters in a latent block model

different, and both of which have the row sizes of at least n_{\min}/K_0 , where n_{\min} is the minimum row size of a block in the **null** block structure. By the same reasoning, for all the **null** column cluster indices $h \in \{1, \dots, H\}$, at least one **estimated** column cluster contains p_h/H_0 or more columns that are assigned to the h th column cluster in the **null** block structure, where p_h is the number of columns in the h th **null** column cluster. By combining these facts, there exists at least one **estimated** block that contains two or more submatrices, both of which have the sizes of at least $(n_{\min}/K_0) \times (p_{\min}/H_0)$ and whose **null** blocks are mutually different.

Let X_1 and X_2 be such submatrices, whose **null** block-wise mean are q_1 and q_2 , respectively. We can assume $q_1 > q_2$ without loss of generality. In matrix \bar{P} , which has the **estimated** block structure, both of X_1 and X_2 have the same values \bar{q} . Here, $|q_2 - \bar{q}| \geq \frac{|q_1 - q_2|}{2}$ holds if $\bar{q} \geq \frac{q_1 + q_2}{2}$, and otherwise $|q_1 - \bar{q}| \geq \frac{|q_1 - q_2|}{2}$. Therefore, for any \bar{q} , there exists at least one submatrix \bar{X} (which is either X_1 or X_2) with a size of at least $(n_{\min}/K_0) \times (p_{\min}/H_0)$, where all the entries are q (which is either q_1 or q_2) in matrix P and

$$|q - \bar{q}| \geq \frac{\min\{\Delta B^{(1)}, \Delta B^{(2)}\}}{2}. \quad (3.28)$$

Here, we used the definitions of $\Delta B^{(1)}$ and $\Delta B^{(2)}$ in (3.3).

Let (k_1, h_1) be the row and column cluster indices of the **estimated** block which contains the above submatrix \bar{X} . We denote the row and column sizes of the (k_1, h_1) th **estimated** block as n_1 and p_1 , respectively. Let $\underline{A}^{(k_1, h_1)}$, $\underline{P}^{(k_1, h_1)}$, $\bar{\underline{P}}^{(k_1, h_1)}$, and $\hat{\underline{P}}^{(k_1, h_1)}$, respectively, be the (k_1, h_1) th **estimated** block of A , P , \bar{P} , and \hat{P} . We define $\hat{q} \equiv \hat{B}_{k_1 h_1}$. In regard to the difference between matrices \bar{P} and \hat{P} (both of which have the **estimated** block structure), we have

$$\begin{aligned} |\hat{q} - \bar{q}| &= \frac{1}{n_1 p_1} \left| \sum_{i=1}^{n_1} \sum_{j=1}^{p_1} \left(\hat{\underline{P}}_{ij}^{(k_1, h_1)} - \bar{\underline{P}}_{ij}^{(k_1, h_1)} \right) \right| \\ &= \frac{1}{n_1 p_1} \left| \sum_{i=1}^{n_1} \sum_{j=1}^{p_1} \left(\underline{A}_{ij}^{(k_1, h_1)} - \underline{P}_{ij}^{(k_1, h_1)} \right) \right| = \frac{1}{n_1 p_1} \left| \left\langle \mathbf{u}_1, \left(\underline{A}^{(k_1, h_1)} - \underline{P}^{(k_1, h_1)} \right) \mathbf{u}_2 \right\rangle \right| \\ &\leq \frac{1}{n_1 p_1} \|\mathbf{u}_1\| \|\mathbf{u}_2\| \|\underline{A}^{(k_1, h_1)} - \underline{P}^{(k_1, h_1)}\|_{\text{op}} = \frac{1}{\sqrt{n_1 p_1}} \|\underline{A}^{(k_1, h_1)} - \underline{P}^{(k_1, h_1)}\|_{\text{op}} \\ &\leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \|\underline{A}^{(k_1, h_1)} - \underline{P}^{(k_1, h_1)}\|_{\text{op}} \leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \|A - P\|_{\text{op}} \\ &\leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \sum_{k=1}^K \sum_{h=1}^H \|A^{(k, h)} - P^{(k, h)}\|_{\text{op}} = \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \sum_{k=1}^K \sum_{h=1}^H S_{kh} \|Z^{(k, h)}\|_{\text{op}} \end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} K H \max_{k=1, \dots, K, h=1, \dots, H} S_{kh} \|Z\|_{\text{op}} = O_p\left(\frac{1}{\sqrt{m}}\right), \quad (3.29)$$

where $A^{(k,h)}$, $P^{(k,h)}$, and $Z^{(k,h)}$, respectively, are the (k, h) th **null** blocks of matrices A , P , and Z , and $\mathbf{u}_1 = [1, 1, \dots, 1]^\top \in \mathbb{R}^{n_1}$ and $\mathbf{u}_2 = [1, 1, \dots, 1]^\top \in \mathbb{R}^{p_1}$. To derive the final equation in (3.29), we used the assumptions that $n_{\min}, p_{\min} = \Omega_p(m)$ and that K, H, K_0 , and H_0 are fixed constants, and the fact that $\|Z\|_{\text{op}} = O_p(\sqrt{m})$ holds from Theorem 2.1.1.

Let $\mathcal{E}_{m,C}$ be the event that $|q - \bar{q}| - C/\sqrt{m} \leq |q - \hat{q}|$ holds. For all q, \bar{q} , and \hat{q} , the following inequality holds:

$$\left| |q - \bar{q}| - |q - \hat{q}| \right| \leq |\hat{q} - \bar{q}|. \quad (3.30)$$

By combining (3.29) and (3.30), for all $\epsilon > 0$, there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $\Pr(\mathcal{E}_{m,C}) \geq 1 - \epsilon$ holds.

From now on, we denote the row and column sizes of submatrix \bar{X} , respectively, by \bar{n}_1 and \bar{p}_1 . Let $A^*, P^*, \bar{P}^*, \hat{P}^*, Z^*$, and \hat{Z}^* , respectively, be the submatrices of matrices $A, P, \bar{P}, \hat{P}, Z$, and \hat{Z} with the same row and column indices as submatrix \bar{X} . We also denote the constant entries of the submatrices of σ and $\hat{\sigma}$ with the same row and column indices as submatrix \bar{X} , respectively, as σ^* and $\hat{\sigma}^*$. From the definition (3.11) and since the operator norm of a submatrix is not larger than that of the original matrix, we have

$$\begin{aligned} \|\hat{Z}\|_{\text{op}} &\geq \|\hat{Z}^*\|_{\text{op}} = \frac{1}{\hat{\sigma}^*} \|A^* - \hat{P}^*\|_{\text{op}} = \frac{1}{\hat{\sigma}^*} \|(A^* - P^*) + (P^* - \hat{P}^*)\|_{\text{op}} \\ &\geq \frac{1}{\hat{\sigma}^*} \left| \|A^* - P^*\|_{\text{op}} - \|P^* - \hat{P}^*\|_{\text{op}} \right| \\ &= \frac{1}{\hat{\sigma}^*} \left| \sigma^* \|Z^*\|_{\text{op}} - \|P^* - \hat{P}^*\|_{\text{op}} \right|. \end{aligned} \quad (3.31)$$

First, the order of the estimated standard deviation $\hat{\sigma}^*$ is given by

$$\hat{\sigma}^* = O_p(1). \quad (3.32)$$

The proof of (3.32) is in Appendix 3.C.

The only non-zero (and thus, the largest) singular value of matrix $(P^* - \hat{P}^*)$ is $\sqrt{\bar{n}_1 \bar{p}_1} |q - \hat{q}|$. Since the largest singular value of a matrix is equal to its operator norm, we have

$$\|P^* - \hat{P}^*\|_{\text{op}} = \sqrt{\bar{n}_1 \bar{p}_1} |q - \hat{q}| \geq \sqrt{\frac{n_{\min} p_{\min}}{K_0 H_0}} |q - \hat{q}|. \quad (3.33)$$

Therefore, by combining this fact with (3.28), if the statement of event $\mathcal{E}_{m,C}$ holds, the following inequality also holds:

$$\sqrt{\frac{n_{\min} p_{\min}}{K_0 H_0}} \left(\frac{\min\{\Delta B^{(1)}, \Delta B^{(2)}\}}{2} - \frac{C}{\sqrt{m}} \right) \leq \|P^* - \hat{P}^*\|_{\text{op}}, \quad (3.34)$$

3. Statistical test on the number of biclusters in a latent block model

which results in that $\|P^* - \hat{P}^*\|_{\text{op}} = \Omega_p(m)$ from the assumption (iv).

Also, from Theorem 2.1.1, we have $\|Z^*\|_{\text{op}} \leq \|Z\|_{\text{op}} = O_p(\sqrt{m})$. By substituting this fact, (3.32), and (3.34) into (3.31), we finally obtain

$$\|\hat{Z}\|_{\text{op}}^2 = \Omega_p(m^2). \quad (3.35)$$

Here, $\|\hat{Z}\|_{\text{op}}^2$ is equal to the maximum eigenvalue $\hat{\lambda}_1$ of $\hat{Z}^\top \hat{Z}$, and the test statistic is $T = \frac{\hat{\lambda}_1 - a^{\text{TW}}}{b^{\text{TW}}}$. Using the definition (2.3), we obtain $a^{\text{TW}} = O_p(m)$ and

$$\begin{aligned} b^{\text{TW}} &= (\sqrt{n} + \sqrt{p}) \left(\frac{1}{\sqrt{n}} + \frac{1}{\sqrt{p}} \right)^{\frac{1}{3}} = (\sqrt{\beta_1 m} + \sqrt{\beta_2 m}) \left(\frac{1}{\sqrt{\beta_1 m}} + \frac{1}{\sqrt{\beta_2 m}} \right)^{\frac{1}{3}} \\ &= m^{\frac{1}{3}} \left(\sqrt{\beta_1} + \sqrt{\beta_2} \right) \left(\frac{1}{\sqrt{\beta_1}} + \frac{1}{\sqrt{\beta_2}} \right)^{\frac{1}{3}}, \end{aligned} \quad (3.36)$$

where we used the definitions $\beta_1 \equiv n/m$ and $\beta_2 \equiv p/m$.

By combining these results and (3.35), we obtain

$$T m^{\frac{1}{3}} = \Omega_p(m^2) \iff T = \Omega_p\left(m^{\frac{5}{3}}\right), \quad (3.37)$$

which concludes the proof. \square

Theorem 3.4.3 (Unrealizable case, upper bound). *Under the assumptions in Section 3.2, if $K_0 < K$ or $H_0 < H$,*

$$T = O_p\left(m^{\frac{5}{3}}\right), \quad (3.38)$$

where T is defined as in (3.12).

Proof. We define P , \bar{P} , and \hat{P} as in Theorem 3.4.2. Let $\hat{Z}^{(k,h)}$, $\underline{A}^{(k,h)}$, and $\hat{P}^{(k,h)}$, respectively, be the (k, h) th **estimated** blocks of matrices \hat{Z} , A , and \hat{P} . We denote the row and column sizes of the (k, h) th **estimated** block as \underline{n}_k and \underline{p}_h , respectively. Since the operator norm of a matrix is not larger than the sum of the operator norms of all its blocks, we have

$$\begin{aligned} \|\hat{Z}\|_{\text{op}} &\leq \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \|\hat{Z}^{(k,h)}\|_{\text{op}} = \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \frac{1}{\hat{S}_{kh}} \|\underline{A}^{(k,h)} - \hat{P}^{(k,h)}\|_{\text{op}} \\ &= \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \frac{\sqrt{\underline{n}_k \underline{p}_h}}{\|\underline{A}^{(k,h)} - \hat{P}^{(k,h)}\|_{\text{F}}} \|\underline{A}^{(k,h)} - \hat{P}^{(k,h)}\|_{\text{op}} \\ &\leq \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \frac{\sqrt{\underline{n}_k \underline{p}_h}}{\|\underline{A}^{(k,h)} - \hat{P}^{(k,h)}\|_{\text{F}}} \|\underline{A}^{(k,h)} - \hat{P}^{(k,h)}\|_{\text{F}} \end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$= \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \sqrt{n_k p_h} \leq K_0 H_0 \sqrt{np} = O_p(m). \quad (3.39)$$

The test statistic is $T = \frac{\hat{\lambda}_1 - a^{\text{TW}}}{b^{\text{TW}}}$, where $\hat{\lambda}_1 = \|\hat{Z}\|_{\text{op}}^2 = O_p(m^2)$. Based on the same discussion as in Theorem 3.4.2, $a^{\text{TW}} = O_p(m)$ and (3.36) hold. Consequently, we obtain $T = O_p(m^2/m^{\frac{1}{3}}) = O_p(m^{\frac{5}{3}})$, which concludes the proof. \square

3.5 Experiments

3.5.1 Realizable case: Convergence of test statistic T in law to Tracy-Widom distribution

First of all, we check the convergence of the proposed test statistic T in law to the TW_1 distribution, under the *realizable* setting, which has been stated in Theorem 3.4.1, by using synthetic data that were generated based on three types of distributions:

- **Gaussian LBM:** The observed matrices were generated whose entries in the (k, h) th block follow the normal distribution $\mathcal{N}(B_{kh}, S_{kh})$. In the Gaussian LBM setting, we used the following null model and parameters:

$$(K, H) = (4, 3), \quad B = \begin{pmatrix} 0.9 & 0.1 & 0.4 \\ 0.2 & 0.7 & 0.3 \\ 0.3 & 0.2 & 0.8 \\ 0.6 & 0.9 & 0.1 \end{pmatrix}, \quad S = \begin{pmatrix} 0.08 & 0.06 & 0.15 \\ 0.14 & 0.12 & 0.07 \\ 0.09 & 0.1 & 0.11 \\ 0.16 & 0.13 & 0.05 \end{pmatrix}. \quad (3.40)$$

- **Bernoulli LBM:** The observed matrices were generated whose entries in the (k, h) th block follow the Bernoulli distribution $\text{Bernoulli}(B_{kh})$. In the Bernoulli LBM setting, we used the following null model and parameters:

$$(K, H) = (4, 3), \quad B = \begin{pmatrix} 0.9 & 0.1 & 0.4 \\ 0.2 & 0.7 & 0.3 \\ 0.3 & 0.2 & 0.8 \\ 0.6 & 0.9 & 0.1 \end{pmatrix}. \quad (3.41)$$

- **Poisson LBM:** The observed matrices were generated whose entries in the (k, h) th block follow the Poisson distribution $\text{Pois}(B_{kh})$. In the Poisson LBM setting, we used the following null model and parameters:

$$(K, H) = (4, 3), \quad B = \begin{pmatrix} 9.0 & 1.0 & 4.0 \\ 2.0 & 7.0 & 3.0 \\ 3.0 & 2.0 & 8.0 \\ 6.0 & 9.0 & 1.0 \end{pmatrix}. \quad (3.42)$$

3. Statistical test on the number of biclusters in a latent block model

Based on the above LBMs, we independently generated 1000 observed matrices, estimated their block structures based on the Ward's hierarchical clustering algorithm [150], and computed the test statistic T . With respect to the matrix size, we tried the following 10 settings: $(n, p) = (300 \times i, 225 \times i)$, $i = 1, \dots, 10$. When generating an observed matrix, the null cluster of each row was independently chosen from the discrete uniform distribution on $\{1, 2, 3, 4\}$. Similarly, the null cluster of each column was independently chosen from the discrete uniform distribution on $\{1, 2, 3\}$.

Figures 3.3, 3.4, and 3.5, respectively, show the Q-Q plots of the test statistic T and the TW_1 distribution in the Gaussian, Bernoulli, and Poisson settings. Each plotted point corresponds to a sample of test statistic T , and the horizontal and vertical lines, respectively, show its theoretical and sample quantiles. These figures show that the test statistic converged in law to the TW_1 distribution.

Figure 3.6 shows the ratios of the trials where $T \geq t(0.01)$, $T \geq t(0.05)$, and $T \geq t(0.1)$ for the above three LBM settings, where $t(\alpha)$ is the α upper quantile of the TW_1 distribution. We used the approximated values $t(0.01) \approx 2.02345$, $t(0.05) \approx 0.97931$, and $t(0.1) \approx 0.45014$, according to Table 2 in [145]. From Figure 3.6, we see that the tail probabilities of the test statistic T also converged to those of the TW_1 distributions for all of the three LBM settings.

We also plotted the results of the Kolmogorov-Smirnov test [36] for the test statistic T in Figure 3.7. We tested whether the distribution of T is the TW_1 distribution or not based on the test statistic $D\sqrt{r}$, where D is the maximum absolute difference between the empirical distribution function of T and the cumulative distribution function of the TW_1 distribution, and r is the sample size, which is set at 1000 in this experiment. Figure 3.7 shows the convergence of the proposed test statistic T in law to the TW_1 distribution under the realizable setting.

3.5.2 Unrealizable case: Asymptotic behavior of test statistic T

Next, we checked the asymptotic behavior of the proposed test statistic T under the *unrealizable* setting, which has been stated in Theorems 3.4.2 and 3.4.3, by using synthetic data that were generated based on the same three types of distributions as in Section 3.5.1. By combining Theorems 3.4.2 and 3.4.3, we obtain the following theorem:

Theorem 3.5.1 (Unrealizable case, two-sided bound). *Under the assumptions in Section 3.2, if $K_0 < K$ or $H_0 < H$,*

$$T = \Theta_p \left(m^{\frac{5}{3}} \right). \quad (3.43)$$

In other words, with high probability, the proposed test statistic T increases in proportion to $m^{\frac{5}{3}}$ in the limit of $m \rightarrow \infty$.

3. Statistical test on the number of biclusters in a latent block model

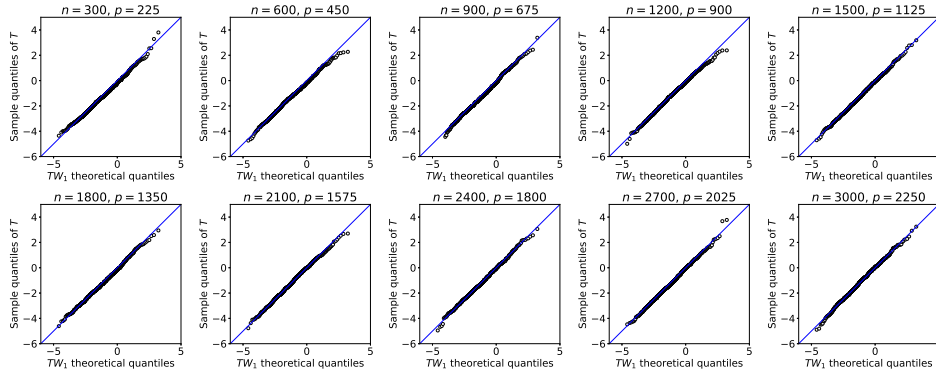


Figure 3.3: Q-Q plot of test statistic T against the TW_1 distribution in the setting of **Gaussian case**.

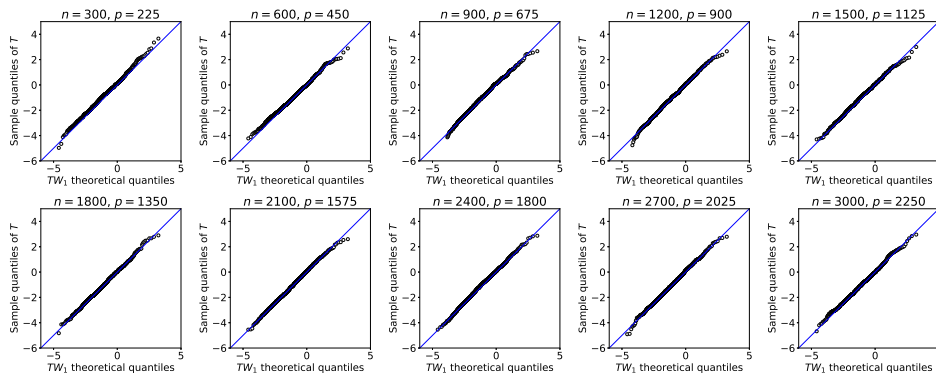


Figure 3.4: Q-Q plot of test statistic T against the TW_1 distribution in the setting of **Bernoulli case**.

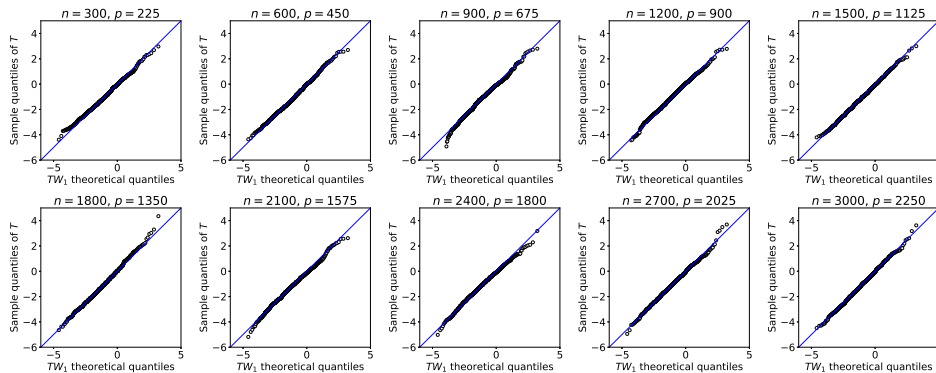


Figure 3.5: Q-Q plot of test statistic T against the TW_1 distribution in the setting of **Poisson case**.

3. Statistical test on the number of biclusters in a latent block model

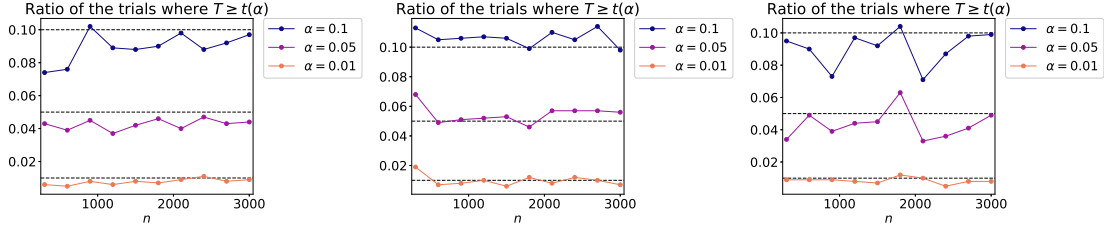


Figure 3.6: Ratio of the number of trials where test statistic $T \geq t(\alpha)$, where $t(\alpha)$ is the α upper quantile of the TW_1 distribution. The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions. The horizontal line shows the number of rows n in the observed matrix.

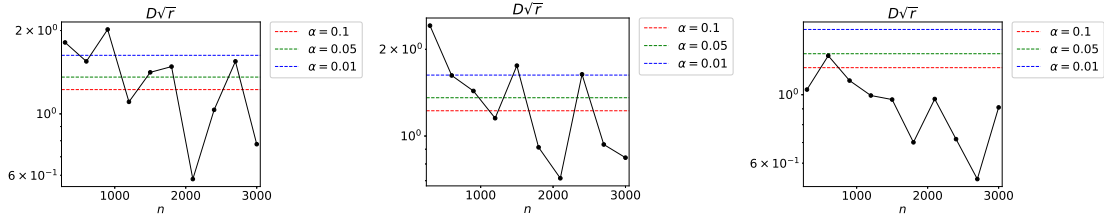


Figure 3.7: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the test statistic T . The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions. If the test statistic $D\sqrt{r}$ is larger than the significance level α , then the null hypothesis that T follows the TW_1 distribution is rejected, and otherwise, the null hypothesis is accepted.

With respect to the null models and parameters, we used the same settings as in Section 3.5.1 for all of the three LBM settings (i.e., Gaussian, Bernoulli, and Poisson LBMs). Based on such settings, we independently generated 100 observed matrices, estimated their block structures based on the Ward's hierarchical clustering algorithm [150], and computed the test statistic T . With respect to the matrix size, we tried the following 10 settings: $(n, p) = (200 \times i, 150 \times i)$, $i = 1, \dots, 10$. When generating an observed matrix, the null cluster of each row was independently chosen from the discrete uniform distribution on $\{1, 2, 3, 4\}$. Similarly, the null cluster of each column was independently chosen from the discrete uniform distribution on $\{1, 2, 3\}$.

Figures 3.8 and 3.9 show the asymptotic behavior of the proposed test statistic T under the unrealizable setting. As shown in Theorem 3.5.1, we see that T increases in proportion to $m^{\frac{5}{3}}$, where $n, p \propto m$.

3. Statistical test on the number of biclusters in a latent block model

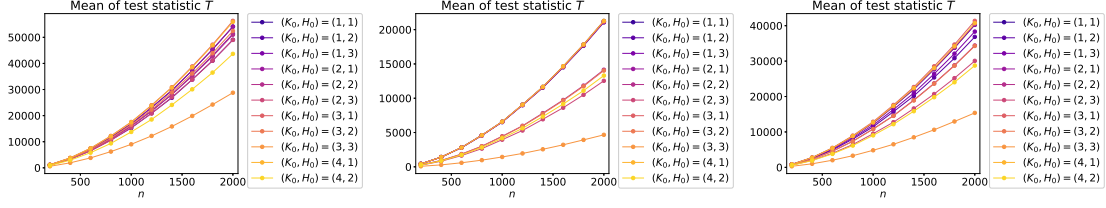


Figure 3.8: Mean test statistic T in the unrealizable case for 100 trials. The null row and column cluster numbers are 4 and 3, respectively. The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions. The horizontal line shows the number of rows n in the observed matrix.

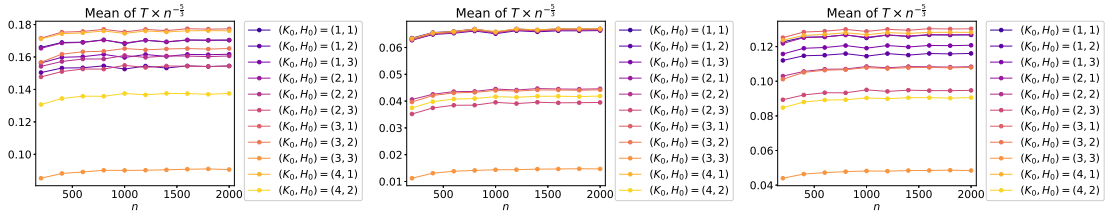


Figure 3.9: Mean test statistic T divided by $n^{\frac{5}{3}}$ in the unrealizable case for 100 trials. The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions.

3.5.3 Accuracy of the proposed goodness-of-fit test

Finally, we evaluated the proposed goodness-of-fit test in terms of its accuracy. By using synthetic data that were generated based on the same three types of distributions as in Section 3.5.1, we checked the ratio of trials where the selected set of cluster numbers (K_0, H_0) is equal to the null one (K, H) . Here, we set the null set of cluster numbers at $(K, H) = (4, 3)$. For each LBM setting (i.e., Gaussian, Bernoulli, and Poisson LBMs), we tried 10 settings with respect to the block-wise mean B . The concrete settings were as follows:

- **Gaussian LBM:** We used the following parameters:

$$B' = \begin{pmatrix} 0.9 & 0.1 & 0.4 \\ 0.2 & 0.7 & 0.3 \\ 0.3 & 0.2 & 0.8 \\ 0.6 & 0.9 & 0.1 \end{pmatrix}, \quad S = \begin{pmatrix} 0.08 & 0.06 & 0.15 \\ 0.14 & 0.12 & 0.07 \\ 0.09 & 0.1 & 0.11 \\ 0.16 & 0.13 & 0.05 \end{pmatrix},$$

$$\forall k, h, B_{kh} = \left(1 - \frac{t}{10}\right) (B'_{kh} - 0.5) + 0.5, \quad \text{for } t = 0, \dots, 9. \quad (3.44)$$

3. Statistical test on the number of biclusters in a latent block model

- **Bernoulli LBM:** We used the following parameters:

$$B' = \begin{pmatrix} 0.9 & 0.1 & 0.4 \\ 0.2 & 0.7 & 0.3 \\ 0.3 & 0.2 & 0.8 \\ 0.6 & 0.9 & 0.1 \end{pmatrix},$$

$$\forall k, h, B_{kh} = \left(1 - \frac{t}{10}\right) (B'_{kh} - 0.5) + 0.5, \quad \text{for } t = 0, \dots, 9. \quad (3.45)$$

- **Poisson LBM:** We used the following parameters:

$$B' = \begin{pmatrix} 9.0 & 1.0 & 4.0 \\ 2.0 & 7.0 & 3.0 \\ 3.0 & 2.0 & 8.0 \\ 6.0 & 9.0 & 1.0 \end{pmatrix},$$

$$\forall k, h, B_{kh} = \left(1 - \frac{t}{10}\right) (B'_{kh} - 5) + 5, \quad \text{for } t = 0, \dots, 9. \quad (3.46)$$

With respect to the matrix size, we tried the following 10 settings for each LBM setting and for each setting of B : $(n, p) = (40 \times i, 30 \times i)$, $i = 1, \dots, 10$. When generating an observed matrix, the null cluster of each row was independently chosen from the discrete uniform distribution on $\{1, 2, 3, 4\}$. Similarly, the null cluster of each column was independently chosen from the discrete uniform distribution on $\{1, 2, 3\}$. In each of 3 (Gaussian, Bernoulli, or Poisson LBM) $\times 10$ (for the setting of B) $\times 10$ (for the setting of matrix size) settings, we generated 1000 observed matrices and applied the proposed sequential goodness-of-fit test, until the null hypothesis $(K, H) = (K_0, H_0)$ was accepted. For each observed matrix, we estimated its block structure based on the Ward's hierarchical clustering algorithm [150] under each setting of a hypothetical set of cluster numbers (K_0, H_0) , computed the test statistic T , and performed the proposed test for the given cluster numbers (K_0, H_0) using a significance level of $\alpha = 0.01$. Figures 3.10, 3.11, and 3.12, respectively, show the examples of generated observed matrices of Gaussian, Bernoulli, and Poisson LBMs.

Figure 6 shows the accuracy of the proposed test under 10 different settings of block-wise mean B . From Figure 6, we see that the test accuracy increases with matrix size n for a fixed block-wise mean B , and that it decreases with the smaller differences between the block-wise means for a fixed matrix size n .

Comparison to the integrated completed likelihood (ICL) We also checked the difference in the behavior of the proposed test and the ICL. For the Bernoulli LBM, we can compute the asymptotic ICL [76] by assuming the following model:

$$p(A, g^{(1)}, g^{(2)} | \pi, \rho, B)$$

3. Statistical test on the number of biclusters in a latent block model

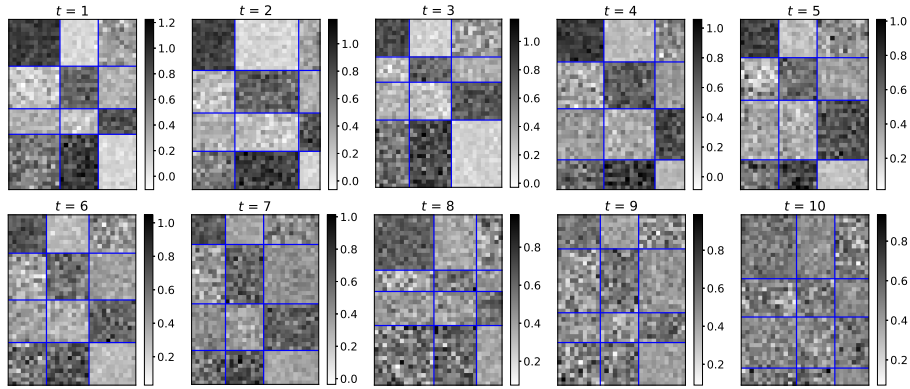


Figure 3.10: Examples of null block structures of the **Gaussian LBM**. 40×30 observed matrices are plotted for 10 different settings of B ($t = 1, \dots, 10$). The rows and columns of matrix A were sorted according to the null clusters.

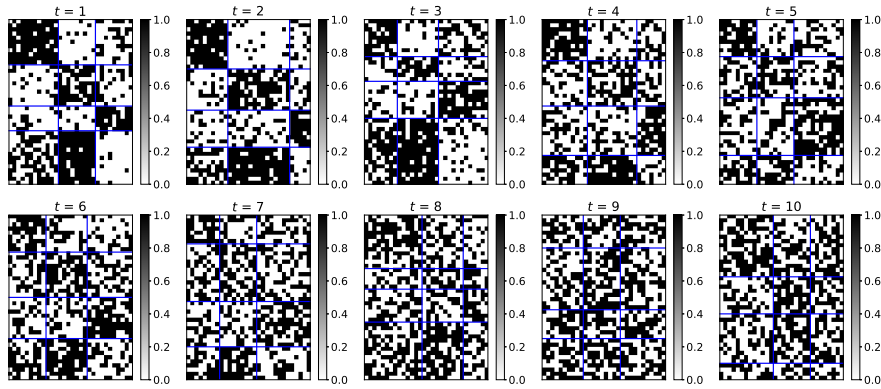


Figure 3.11: Examples of null block structures of the **Bernoulli LBM**.

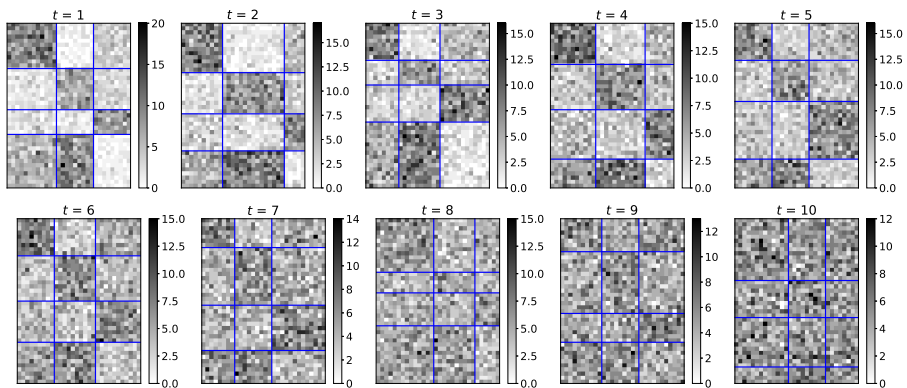


Figure 3.12: Examples of null block structures of the **Poisson LBM**.

3. Statistical test on the number of biclusters in a latent block model

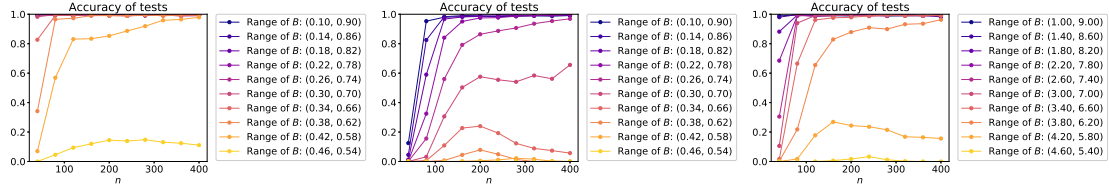


Figure 3.13: Accuracy of the proposed goodness-of-fit test under 10 different settings of block-wise mean B . The left, center, and right figures, respectively, show the results for the settings of Gaussian, Bernoulli, and Poisson distributions.

$$\begin{aligned}
 &= \left(\prod_k \pi_k^{|I_k|} \right) \left(\prod_h \rho_h^{|J_h|} \right) \left[\prod_{i,j} B_{g_i^{(1)} g_j^{(2)}}^{A_{ij}} \left(1 - B_{g_i^{(1)} g_j^{(2)}} \right)^{1-A_{ij}} \right], \\
 p(\pi) &= \frac{\Gamma(a_D K_0)}{\Gamma(a_D)^{K_0}} \prod_{k=1}^{K_0} \pi_k^{a_D - 1}, \quad p(\rho) = \frac{\Gamma(a_D H_0)}{\Gamma(a_D)^{H_0}} \prod_{h=1}^{H_0} \rho_h^{a_D - 1}, \\
 p(B) &= \prod_{k,h} \frac{B_{kh}^{b_B - 1} (1 - B_{kh})^{b_B - 1}}{B(b_B, b_B)}, \tag{3.47}
 \end{aligned}$$

where $p(\cdot)$ represents a probability density, and a_D and b_B are the hyperparameters.

From Lemma 4.2 in [76], for an estimated block structure $(\hat{g}^{(1)}, \hat{g}^{(2)})$, the resulting asymptotic ICL is given by

$$\begin{aligned}
 \text{ICL}(K_0, H_0) &= \sum_k |I_k| \log \left(\frac{|I_k|}{n} \right) + \sum_h |J_h| \log \left(\frac{|J_h|}{p} \right) \\
 &\quad + \sum_{k,h} |I_k| |J_h| \left[\hat{B}_{kh} \log \hat{B}_{kh} + (1 - \hat{B}_{kh}) \log (1 - \hat{B}_{kh}) \right] \\
 &\quad - \frac{K_0 - 1}{2} \log n - \frac{H_0 - 1}{2} \log p - \frac{K_0 H_0}{2} \log(np). \tag{3.48}
 \end{aligned}$$

The proof of (3.48) is given in Appendix 3.D.

As a preliminary experiment, we checked the relationship between the proposed test statistic T or the ICL and the eigenvalues of matrix $\hat{Z}^\top \hat{Z}$. By setting the cluster numbers and the matrix size at $(K, H) = (K_0, H_0) = (4, 3)$ and $(n, p) = (300, 225)$, respectively, we generated 1000 observed matrices based on the Bernoulli distribution and estimated their bicluster structure as in Section 3.5.1. For each estimated bicluster structure, we computed the proposed test statistic T , the ICL, and the eigenvalues of matrix $\hat{Z}^\top \hat{Z}$. Figure 3.14 shows the correlation coefficients ρ between the proposed test statistic T or the ICL and the eigenvalues of matrix $\hat{Z}^\top \hat{Z}$. Since the proposed test statistic T is a scaled maximum eigenvalue of $\hat{Z}^\top \hat{Z}$, the correlation coefficient between T and the maximum eigenvalue is always one, while those between T and the other eigenvalues were relatively

3. Statistical test on the number of biclusters in a latent block model

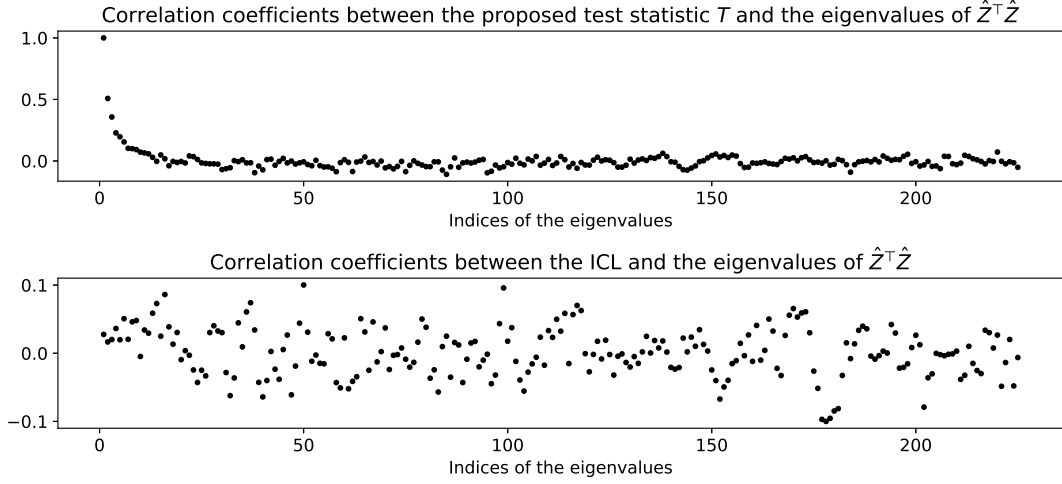


Figure 3.14: Correlation coefficients between the proposed test statistic T or the ICL and the eigenvalues of matrix $\hat{Z}^\top \hat{Z}$. As for the horizontal axes, index $j \in \{1, \dots, p\}$ corresponds to the j th largest eigenvalue.

small. In regard to the ICL, unlike the proposed test statistic, the correlation coefficients were small (i.e., $-0.2 < \rho < 0.2$) for all the indices of eigenvalues.

Next, to check the accuracy of the proposed test and the ICL, we generated synthetic binary data matrices based on the Bernoulli distribution as in Section 3.5.3, and checked the ratio of trials where the selected set of cluster numbers (K_0, H_0) is equal to the null one (K, H) . We set the null set of cluster numbers at $(K, H) = (4, 3)$, and tried the following five settings with respect to the block-wise mean B .

$$B' = \begin{pmatrix} 0.9 & 0.1 & 0.4 \\ 0.2 & 0.7 & 0.3 \\ 0.3 & 0.2 & 0.8 \\ 0.6 & 0.9 & 0.1 \end{pmatrix},$$

$$\forall k, h, B_{kh} = \left(1 - \frac{t}{5}\right) (B'_{kh} - 0.5) + 0.5, \quad \text{for } t = 0, \dots, 4. \quad (3.49)$$

With respect to the matrix size, we tried the following five settings for each setting of B : $(n, p) = (40 \times i, 30 \times i)$, $i = 1, \dots, 5$. The null block of each element was chosen in the same way as in Section 3.5.3. In each of 5 (for the setting of B) \times 5 (for the setting of matrix size) settings, we generated 100 observed matrices, and applied the proposed test using a significance level of $\alpha = 0.01$ and the model selection based on the ICL. Unlike the proposed sequential test, which stopped if the null hypothesis was accepted, the ICL was computed for all the sets of cluster numbers from $(1, 1)$ to (n, p) and then the optimal setting was selected that achieved the maximum ICL. For each setting, we estimated the

3. Statistical test on the number of biclusters in a latent block model

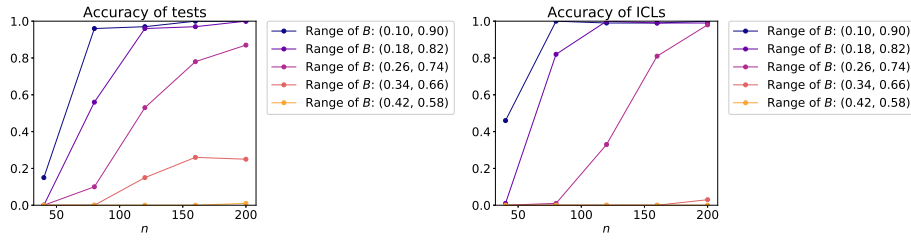


Figure 3.15: Accuracy of the proposed test and the model selection based on the ICL under five different settings of block-wise mean B .

block structure of an observed matrix based on the Ward’s hierarchical clustering algorithm [150].

Figure 3.15 shows the accuracy of the proposed test and the model selection based on the ICL. Although the purpose of the proposed test is not to achieve high accuracy in model selection, in some cases with small differences between the block-wise means $\{B_{kh}\}$, it achieved better performance than the ICL. With larger difference among $\{B_{kh}\}$, the ICL performed better than the proposed test in terms of model selection. Figures 3.16 and 3.17, respectively, show the ratios of the trials where each set of cluster numbers was selected by the proposed test and the ICL. From Figures 3.16 and 3.17, we see that in most cases (e.g., $B_{kh} \in [0.26, 0.74]$ for all (k, h)), the ICL tended to select smaller sets of cluster numbers than the proposed test.

3.5.4 Real data analysis: Congressional Voting Records Data Set

Selection of the number of biclusters by the proposed test and the ICL

We also checked the result when we applied the proposed test to 1984 United States Congressional Voting Records Database from UCI Machine Learning Repository [44]. The original data set contains three types of votes (“yea,” “nay,” and unknown) for the pairs of a congressman and an attribute. We treated unknown as “nay,” as in [159]. The number of instances or congressmen and that of attributes are 435 and 16, respectively. Based on this data set, we defined a binary matrix $A \in \mathbb{R}^{435 \times 16}$, where the elements of one and zero, respectively, correspond to “yea” and “nay.”

As in Section 3.5.3, we applied the proposed sequential tests using a significance level of $\alpha = 0.01$, until the null hypothesis was accepted. We also computed the ICL for each setting of a hypothetical set of cluster numbers (K_0, H_0) , and selected one with the largest ICL. For each setting of a hypothetical set of cluster numbers (K_0, H_0) , we estimated the block structure based on the Ward’s hierarchical clustering algorithm [150].

As a result, the sets of cluster numbers (9, 14) and (3, 13) were selected by the proposed test and the ICL, respectively. Figure 3.18 shows the observed data matrix and its estimated block structures with the selected sets of cluster numbers. From Figure 3.18, we see that

3. Statistical test on the number of biclusters in a latent block model

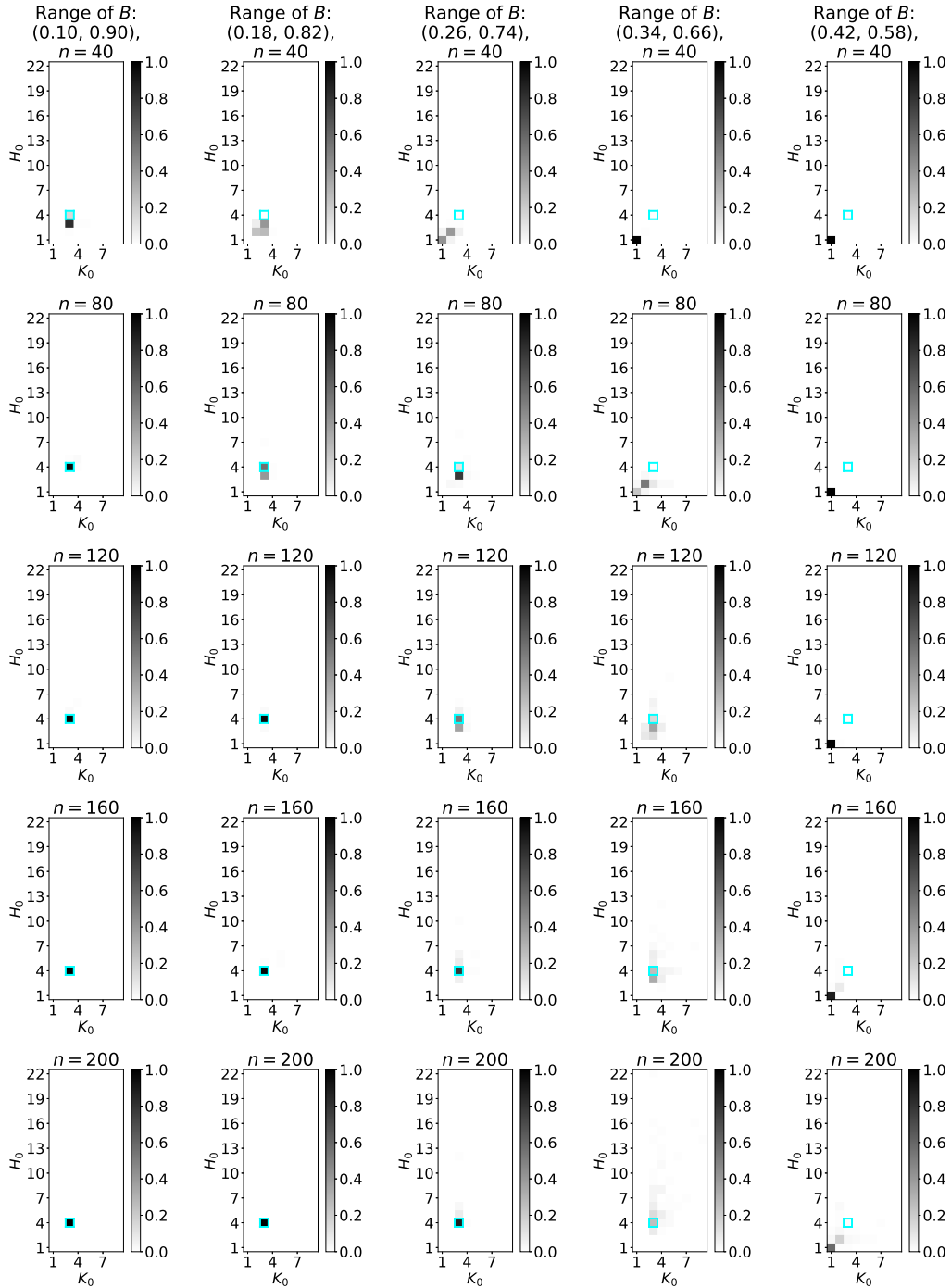


Figure 3.16: Ratios of the trials where each set of cluster numbers (K_0, H_0) was selected by the proposed test. The cyan rectangles show the null set of cluster numbers (K, H) .

3. Statistical test on the number of biclusters in a latent block model

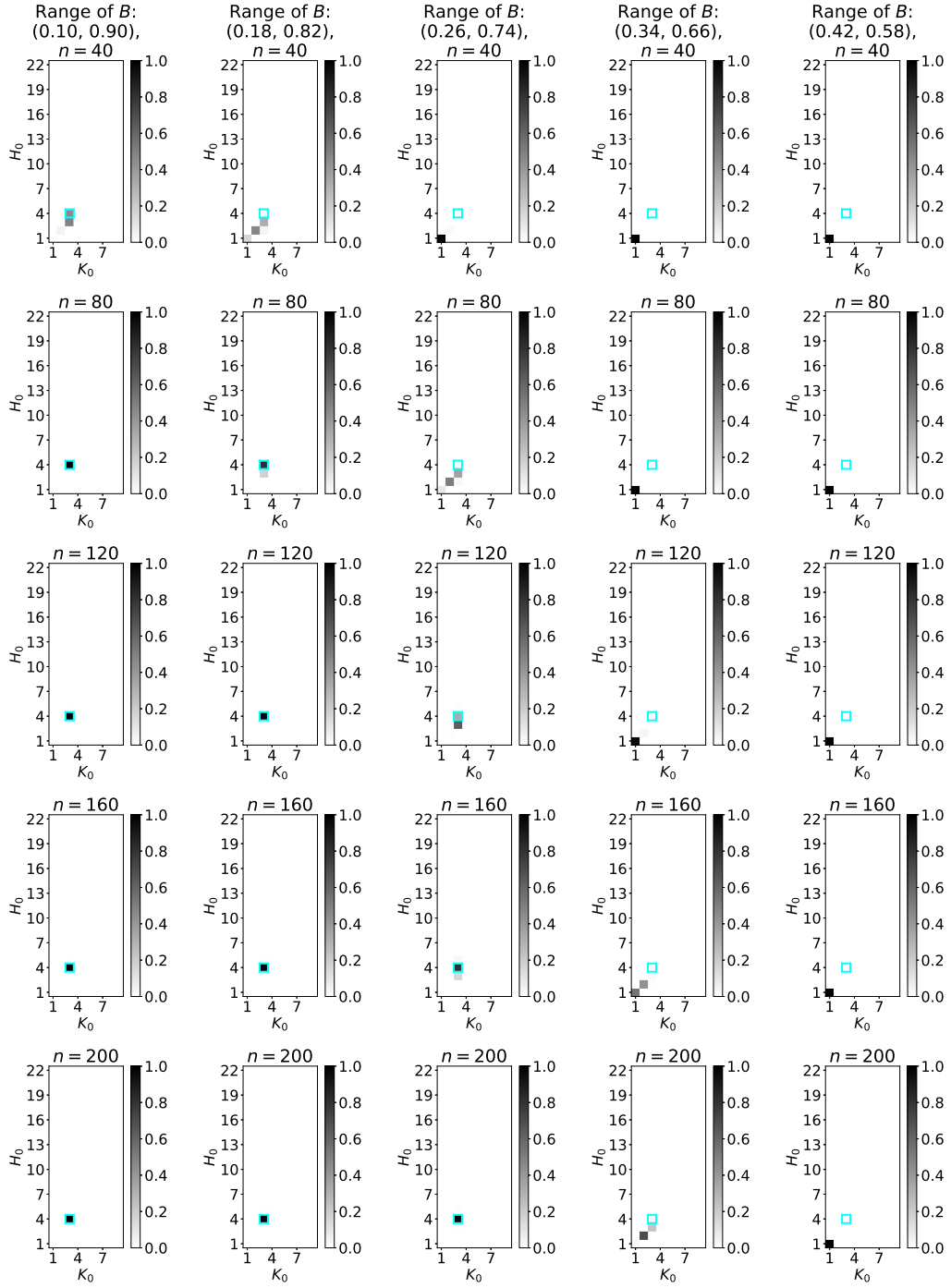


Figure 3.17: Ratios of the trials where each set of cluster numbers (K_0, H_0) was selected by the model selection based on the ICL. The cyan rectangles show the null set of cluster numbers (K, H) .

3. Statistical test on the number of biclusters in a latent block model

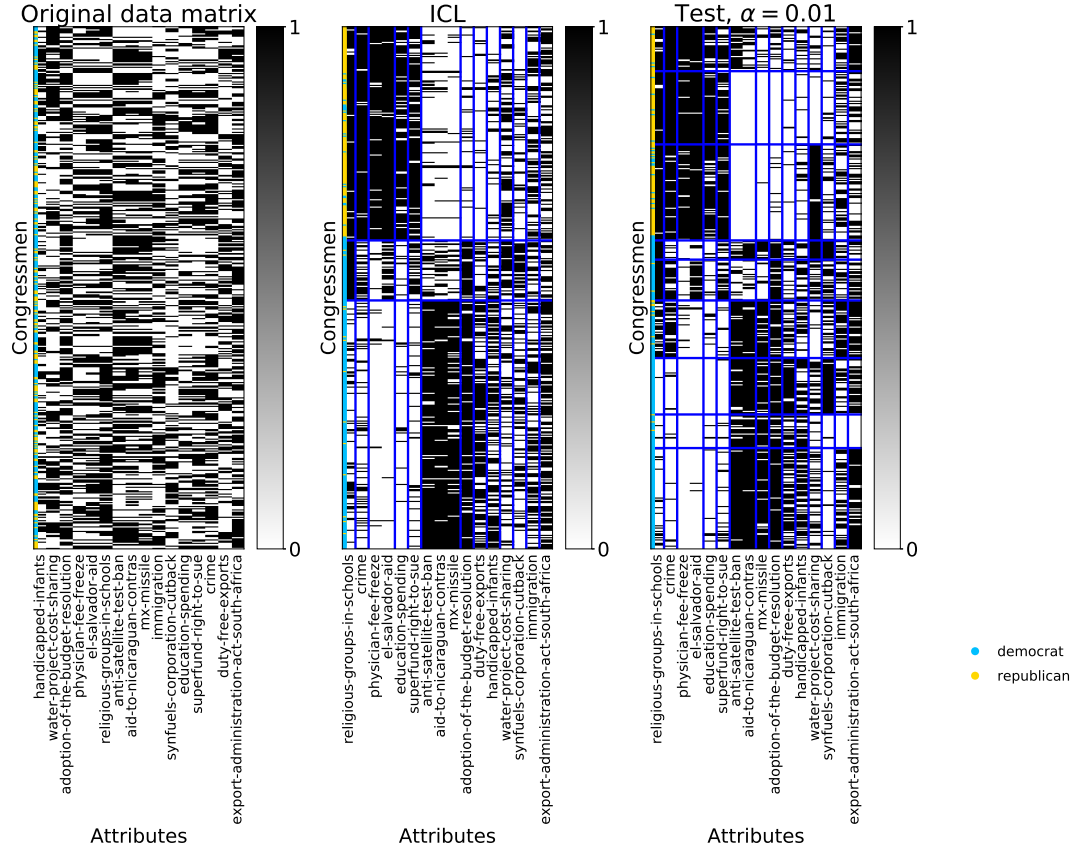


Figure 3.18: The observed data matrix of the Congressional Voting Records Data Set [44] and its estimated block structures. The black and white elements, respectively, show “yea” and “nay.”

a finer block structure was accepted by the proposed test than the ICL, particularly for the row (i.e., congressman) cluster assignments. As for the column (i.e., attribute) cluster assignments, “anti-satellite-test-ban,” “aid-to-nicaraguan-contras,” and “mx-missile” were assigned into the same cluster in the selected block structure of the ICL, whereas the proposed test distinguished the first two attributes from the last one. Figure 3.19 shows the p -value of the proposed test and the ICL for each setting of a hypothetical set of cluster numbers (K_0, H_0) until the null hypothesis was accepted.

Robustness of the proposed test and the ICL to data noise

By using the same Congressional Voting Records Data Set, we compared the proposed test and the ICL in terms of robustness to the data noise. Specifically, we chose \tilde{m} entries without replacement from the discrete uniform distribution on all the entries of matrix A

3. Statistical test on the number of biclusters in a latent block model

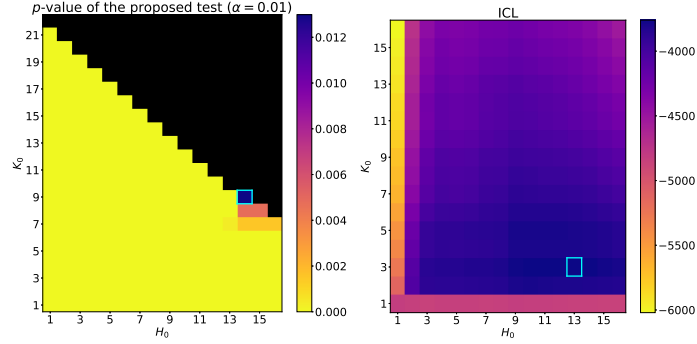


Figure 3.19: The p -value of the proposed test and the ICL for each setting of a hypothetical set of cluster numbers (K_0, H_0) . In the left figure, we plotted the p -values only for the tested settings (i.e., until the null hypothesis was accepted). As for the ICL, we plotted a part of results ($K_0 \leq 16, H_0 \leq 16$) for visibility. The cyan rectangles show the selected sets of cluster numbers.

and flipped the values of these selected entries (i.e., from zero to one or from one to zero). In the experiment, we tried $\check{m} = 2t + 2$ for $t = 0, 1, \dots, 4$. For each setting of \check{m} , we generated 1000 matrices \check{A} with noise based on the above procedure. We estimated the block structure of matrix \check{A} based on the Ward's hierarchical clustering algorithm [150] and applied the proposed test with a significance level of $\alpha = 0.01$ and the ICL.

Figure 3.20 shows the ratio of the trials in which each set of row and column cluster numbers was selected by the proposed test and the ICL. From this result, we see that the ICL was more robust than the proposed test with regard to the data noise. The selected sets of cluster numbers by the ICL were $(3, 9)$ in most trials, while those by the proposed test varied more greatly depending on the selection results of the flipped entries.

3.6 Discussions

In this section, we discuss the proposed test method in terms of the test statistic and the conditions for the generative model.

With respect to the asymptotic behavior, the proposed test has a favorable property in terms of the power. From Theorem 3.5.1, under the alternative hypothesis, the test statistic T increases in proportion to $m^{\frac{5}{3}}$ with high probability, where $n, p \propto m$. In other words, the probability that the test makes a type II error (i.e., $T < t(\alpha)$) converges to zero in the limit of $p \rightarrow \infty$. Based on this fact, in the asymptotic sense, we do not need to consider the correction for the multiple comparison when applying the proposed sequential testing. However, it has not been shown what occurs in the non-asymptotic setting. In general, practical data matrices have finite sizes, where there has been shown no theoretical guarantee like Theorems 3.4.1, 3.4.2, and 3.4.3. On the other hand, for a Gaussian case

3. Statistical test on the number of biclusters in a latent block model

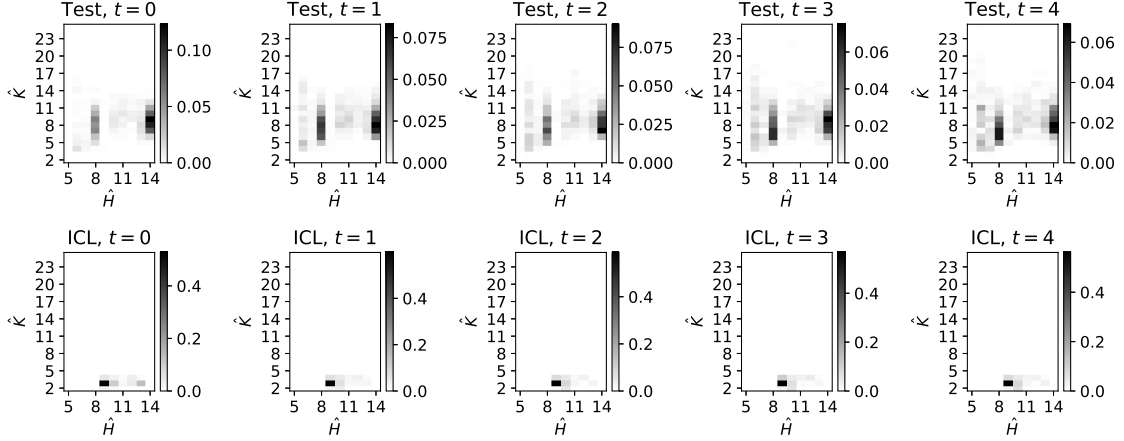


Figure 3.20: The ratio of the trials in which each set of row and column cluster numbers was selected by the proposed test and the ICL. We plotted the results only for the settings which were selected at least once by the proposed test or the ICL.

(i.e., each entry of a matrix independently follows $\mathcal{N}(0, 1)$), the following statement holds [105]: Suppose $n = n(p) > p$ and $n/p \rightarrow \gamma \in [1, \infty)$ in the limit of $p \rightarrow \infty$. Then, for any s_0 , there exists $N_0 \in \mathbb{N}$ such that when $\max(n, p) \geq N_0$ and $\max(n, p)$ is even, for all $s \geq s_0$,

$$|\Pr(T^* \leq s) - F_1(s)| \leq C(s_0)[\max(n, p)]^{-2/3} \exp\left(-\frac{s}{2}\right), \quad (3.50)$$

where T^* is defined as in (2.2), F_1 is the cumulative distribution function of the TW_1 distribution, and $C(\cdot)$ is a continuous and non-increasing function. From the above inequality (3.50), if the clustering algorithm outputs the correct block assignments, the convergence rate of the normalized maximum eigenvalue T^* of matrix $\tilde{Z}^\top \tilde{Z}$ (where \tilde{Z} is defined as in (3.14)) to the TW_1 distribution is $O(m^{-2/3})$. However, since the distribution of T is unknown in the case where the correct block assignment is *not* obtained, the convergence rate of T is also unknown. Deriving the convergence rate of T by considering the above discussion is a future research topic.

In regard to the conditions for using the proposed test method, our proposed test is applicable to a wide range of practical settings (e.g., Bernoulli distribution for binary data matrices and Poisson distribution for sparse ones). Nevertheless, it still requires some assumptions for the latent block structure of an observed matrix. For instance, the row and column cluster numbers (K, H) should be constants that do not increase with the matrix sizes n and p . Also, there should be no too small block (i.e., $n_{\min} = \Omega_p(m)$ and $p_{\min} = \Omega_p(m)$). In some practical cases, where it is more appropriate to assume that the number of blocks increases with the matrix size, it will be useful to construct a test which does not require the above conditions. As for the sub-exponential condition, Ding and Yang

3. Statistical test on the number of biclusters in a latent block model

[43] have shown more relaxed sufficient condition for the scaled maximum eigenvalue T^* to converge in law to the TW_1 distribution. However, the *delocalization property* of an eigenvector of matrix $Z^\top Z$ [17], which we used in Appendix 3.B to prove our main result, has not been derived in the form as in the sub-exponential case [17]. If Theorem 2.1.2 is shown in the above more general case, it would also be possible to extend our proposed test to such a case. Furthermore, there are proposed variants of latent block models with which we assume different block structures from a regular grid [112, 129]. To construct test methods for the above settings is an important topic for future research.

3.7 Chapter conclusion

Latent block models are effective tools for biclustering, where rows and columns of an observed matrix are simultaneously decomposed into clusters. Such a bicluster structure appears in various types of relational data, such as the customer-product transaction data and the document-word relationship data. One open problem in using latent block models is that there has been no statistical test method for determining the number of blocks. In this chapter, we developed a goodness-of-fit test for latent block models based on a result from the random matrix theory. By defining the test statistic T based on the estimators of the block-wise means and standard deviations, we have derived its asymptotic behavior in both realizable (i.e., $(K, H) = (K_0, H_0)$) and unrealizable (i.e., $K > K_0$ or $H > H_0$) cases. Particularly, it has been shown that the test statistic T converges in law to the TW_1 distribution in the realizable case. Based on these results, it was made possible to test whether the given observed matrix had $K_0 \times H_0$ latent blocks or more ones. In the experiments, we showed the validity of the proposed test method in terms of both the asymptotic behavior of the test statistic and the test accuracy by using synthetic data matrices. We also applied the proposed test to practical data and analyzed the result.

3.A Proof of $|\tilde{S}_{kh} - S_{kh}| = O_p\left(\frac{1}{m}\right)$

Let n_k and p_h , respectively, be the row and column sizes of the (k, h) th **null** block, and $A^{(k,h)}$, $P^{(k,h)}$, and $\tilde{P}^{(k,h)}$, respectively, be the (k, h) th **null** blocks of matrices A , P , and \tilde{P} . Here, we prove the following lemma:

Lemma A1. *Under the assumption that the fourth moment of the noise Z_{ij} is bounded ($\mathbb{E}[Z_{ij}^4] < \infty$),*

$$|\tilde{S}_{kh} - S_{kh}| = O_p\left(\frac{1}{m}\right), \quad (3.51)$$

where $\tilde{S}_{kh} = \sqrt{\frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - \tilde{B}_{kh}\right)^2}$.

3. Statistical test on the number of biclusters in a latent block model

Proof. From the above definition of \tilde{S}_{kh} , we have

$$\begin{aligned}\tilde{S}_{kh}^2 &= \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - \tilde{B}_{kh} \right)^2 \\ &= \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - \left(B_{kh} - \tilde{B}_{kh} \right)^2.\end{aligned}\quad (3.52)$$

To derive the second equation, we used the fact that $\tilde{B}_{kh} = \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} A_{ij}^{(k,h)}$. Therefore, the following inequality holds:

$$\left| \tilde{S}_{kh}^2 - S_{kh}^2 \right| \leq \left| \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2 \right| + \left(B_{kh} - \tilde{B}_{kh} \right)^2. \quad (3.53)$$

The first term in (3.53) is given by

$$\begin{aligned}& \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2 \\ &= \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left[\left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2 \right] = \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} Y_{ij}^{(k,h)},\end{aligned}\quad (3.54)$$

where we defined that $Y_{ij}^{(k,h)} \equiv \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2$. Note that $(Y_{ij}^{(k,h)})_{1 \leq i \leq n_k, 1 \leq j \leq p_h}$ is independent. The expectation and the variance of $Y_{ij}^{(k,h)}$ are given by

$$\begin{aligned}\mathbb{E} \left[Y_{ij}^{(k,h)} \right] &= \mathbb{E} \left[\left(A_{ij}^{(k,h)} - B_{kh} \right)^2 \right] - S_{kh}^2 = 0, \\ \mathbb{V} \left[Y_{ij}^{(k,h)} \right] &= \mathbb{E} \left[\left(Y_{ij}^{(k,h)} \right)^2 \right] = \mathbb{E} \left[\left\{ \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2 \right\}^2 \right] \\ &= S_{kh}^4 \left(\mathbb{E} \left[\left(Z_{ij}^{(k,h)} \right)^4 \right] - 1 \right).\end{aligned}\quad (3.55)$$

From (3.55), we have

$$\begin{aligned}\mathbb{E} \left[\frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} Y_{ij}^{(k,h)} \right] &= 0, \\ \mathbb{V} \left[\frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} Y_{ij}^{(k,h)} \right] &= \frac{1}{n_k p_h} S_{kh}^4 \left(\mathbb{E} \left[\left(Z_{ij}^{(k,h)} \right)^4 \right] - 1 \right).\end{aligned}\quad (3.56)$$

3. Statistical test on the number of biclusters in a latent block model

From (3.54), (3.56), and Chebyshev's inequality, for all $t > 0$,

$$\Pr \left(\left| \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2 \right| \geq t \frac{1}{\sqrt{n_k p_h}} \sqrt{S_{kh}^4 \left(\mathbb{E} \left[\left(Z_{ij}^{(k,h)} \right)^4 \right] - 1 \right)} \right) \leq \frac{1}{t^2}. \quad (3.57)$$

Therefore, we have

$$\left| \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(A_{ij}^{(k,h)} - B_{kh} \right)^2 - S_{kh}^2 \right| = O_p \left(\frac{1}{m} \right). \quad (3.58)$$

On the other hand, the second term in (3.53) is given by

$$\left(B_{kh} - \tilde{B}_{kh} \right)^2 = O_p \left(\frac{1}{m^2} \right). \quad (3.59)$$

from (3.16).

By combining (3.53), (3.58), and (3.59),

$$\left| \tilde{S}_{kh}^2 - S_{kh}^2 \right| = O_p \left(\frac{1}{m} \right). \quad (3.60)$$

The difference between \tilde{S}_{kh} and S_{kh} is given by

$$\left| \tilde{S}_{kh} - S_{kh} \right| = \frac{\left| \tilde{S}_{kh}^2 - S_{kh}^2 \right|}{\tilde{S}_{kh} + S_{kh}}. \quad (3.61)$$

Here, from (3.60), $m \left| \tilde{S}_{kh}^2 - S_{kh}^2 \right|$ is bounded in probability. Therefore, \tilde{S}_{kh} converges in probability to S_{kh} , which results in that $\tilde{S}_{kh} + S_{kh} = \Theta_p(1)$. By combining this fact with (3.60) and (3.61), we finally obtain

$$\left| \tilde{S}_{kh} - S_{kh} \right| = O_p \left(\frac{1}{m} \right), \quad (3.62)$$

which concludes the proof. \square

3.B Proof of $\frac{|\lambda_1 - \hat{\lambda}_1|}{b^{\text{TW}}} = O_p \left(m^{-\frac{1}{2\Gamma} + \epsilon} \right)$ for all $\epsilon \in \left(0, \frac{2}{7} \right)$ in realizable case

We first derive the relationship between the maximum eigenvalues of matrices $Z^\top Z$ and $\tilde{Z}^\top \tilde{Z}$ in Lemma B1 and B2.

3. Statistical test on the number of biclusters in a latent block model

Lemma B1. Let λ_1 and $\tilde{\lambda}_1$, respectively, be the maximum eigenvalues of matrices $Z^\top Z$ and $\tilde{Z}^\top \tilde{Z}$ (i.e., $\|Z\|_{\text{op}}^2$ and $\|\tilde{Z}\|_{\text{op}}^2$, respectively). Then, for all $\epsilon \in (0, \frac{1}{2})$, the following equation holds:

$$\lambda_1 \leq \tilde{\lambda}_1 + O_p(m^\epsilon). \quad (3.63)$$

Proof. Let \mathbf{v} and $\tilde{\mathbf{v}}$, respectively, be the normalized eigenvectors of $Z^\top Z$ and $\tilde{Z}^\top \tilde{Z}$, corresponding to the maximum eigenvalues λ_1 and $\tilde{\lambda}_1$:

$$\begin{aligned} Z^\top Z \mathbf{v} &= \lambda_1 \mathbf{v}, \quad \|\mathbf{v}\| = 1, \\ \tilde{Z}^\top \tilde{Z} \tilde{\mathbf{v}} &= \tilde{\lambda}_1 \tilde{\mathbf{v}}, \quad \|\tilde{\mathbf{v}}\| = 1. \end{aligned} \quad (3.64)$$

Since $\sqrt{\tilde{\lambda}_1}$ is the largest singular value of matrix \tilde{Z} , we have

$$\sqrt{\tilde{\lambda}_1} = \sup_{\mathbf{u} \in \mathbb{R}^p} \frac{\|\tilde{Z} \mathbf{u}\|}{\|\mathbf{u}\|} \geq \frac{\|\tilde{Z} \mathbf{v}\|}{\|\mathbf{v}\|} = \|\tilde{Z} \mathbf{v}\| \iff \tilde{\lambda}_1 \geq \|\tilde{Z} \mathbf{v}\|^2. \quad (3.65)$$

We also define the following matrix $Q^{(k,h)}$ for each (k, h) th block:

$$\begin{aligned} Q^{(k,h)} &\equiv Z^{(k,h)} - \frac{\tilde{S}_{kh}}{S_{kh}} \tilde{Z}^{(k,h)} = \frac{1}{S_{kh}} \left(\tilde{P}^{(k,h)} - P^{(k,h)} \right) \\ &= \frac{1}{n_k p_h} \left[\sum_{i \in I_k, j \in J_h} \frac{A_{ij} - P_{ij}}{S_{kh}} \right] \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{bmatrix} \\ &= \left(\frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} Z_{ij}^{(k,h)} \right) \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{bmatrix}, \end{aligned} \quad (3.66)$$

where n_k and p_h , respectively, are the row and column sizes of the (k, h) th **null** block.

Let $\underline{Z}^{(k,h)}$, $\tilde{\underline{Z}}^{(k,h)}$, and $\underline{Q}^{(k,h)}$, respectively be $n \times p$ matrices whose (k, h) th null blocks are $Z^{(k,h)}$, $\tilde{Z}^{(k,h)}$ and $Q^{(k,h)}$ and whose all the other entries are zero. As shown in Figure 3.B1, we define matrix Q as $Q \equiv \sum_{k=1}^K \sum_{h=1}^H \underline{Q}^{(k,h)}$. We also define that $\tau_{kh} \equiv \frac{S_{kh}}{\tilde{S}_{kh}}$.

From (3.65), we have

$$\begin{aligned} \tilde{\lambda}_1 &\geq \|\tilde{Z} \mathbf{v}\|^2 = \left\| \sum_{k=1}^K \sum_{h=1}^H \tilde{\underline{Z}}^{(k,h)} \mathbf{v} \right\|^2 = \left\| \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} (\underline{Z}^{(k,h)} - \underline{Q}^{(k,h)}) \mathbf{v} \right\|^2 \\ &\geq \left[\left\| \sum_{k=1}^K \sum_{h=1}^H (\underline{Z}^{(k,h)} - \underline{Q}^{(k,h)}) \mathbf{v} \right\| - \left\| \sum_{k=1}^K \sum_{h=1}^H (1 - \tau_{kh}) (\underline{Z}^{(k,h)} - \underline{Q}^{(k,h)}) \mathbf{v} \right\| \right]^2 \end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

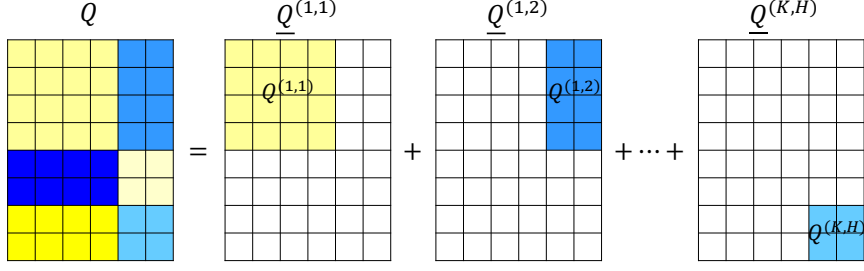


Figure 3.B1: Definition of matrix Q .

$$\begin{aligned}
&= \left[\|(Z - Q)\mathbf{v}\| - \left\| \sum_{k=1}^K \sum_{h=1}^H (1 - \tau_{kh}) (\underline{Z}^{(k,h)} - \underline{Q}^{(k,h)})\mathbf{v} \right\| \right]^2 \\
&\geq \|(Z - Q)\mathbf{v}\|^2 - 2\|(Z - Q)\mathbf{v}\| \left\| \sum_{k=1}^K \sum_{h=1}^H (1 - \tau_{kh}) (\underline{Z}^{(k,h)} - \underline{Q}^{(k,h)})\mathbf{v} \right\| \\
&\geq \|(Z - Q)\mathbf{v}\|^2 - 2\|(Z - Q)\mathbf{v}\| \left[\sum_{k=1}^K \sum_{h=1}^H |1 - \tau_{kh}| \|(\underline{Z}^{(k,h)} - \underline{Q}^{(k,h)})\mathbf{v}\| \right] \\
&\geq \|(Z - Q)\mathbf{v}\|^2 - 2\|(Z - Q)\mathbf{v}\| \left[\sum_{k=1}^K \sum_{h=1}^H |1 - \tau_{kh}| \left(\|\underline{Z}^{(k,h)}\mathbf{v}\| + \|\underline{Q}^{(k,h)}\mathbf{v}\| \right) \right] \\
&\geq \|(Z - Q)\mathbf{v}\|^2 - 2\|(Z - Q)\mathbf{v}\| \left[\sum_{k=1}^K \sum_{h=1}^H |1 - \tau_{kh}| \left(\sqrt{\lambda_1^{(k,h)}} + \|\underline{Q}^{(k,h)}\mathbf{v}\| \right) \right] \\
&\geq \lambda_1 - 2\sqrt{\lambda_1}\|Q\mathbf{v}\| - 2(\sqrt{\lambda_1} + \|Q\mathbf{v}\|) \left[\sum_{k=1}^K \sum_{h=1}^H |1 - \tau_{kh}| \left(\sqrt{\lambda_1^{(k,h)}} + \|\underline{Q}^{(k,h)}\mathbf{v}\| \right) \right], \tag{3.67}
\end{aligned}$$

where $\lambda_1^{(k,h)}$ is the maximum eigenvalue of matrix $(Z^{(k,h)})^\top Z^{(k,h)}$.

From now on, we prove that $\|Q\mathbf{v}\| = O_p\left(\frac{1}{\sqrt{m}}\right)$ and $\|\underline{Q}^{(k,h)}\mathbf{v}\| = O_p\left(\frac{1}{\sqrt{m}}\right)$. We use the following notations:

$$\nu_{kh} \equiv \frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} Z_{ij}^{(k,h)}, \quad \omega_{hh'} \equiv \sum_{k=1}^K n_k \nu_{kh} \nu_{kh'}, \quad \zeta_h \equiv \sum_{h'=1}^H \omega_{hh'} \sum_{j \in J_{h'}} v_j, \tag{3.68}$$

where v_j is the j th entry of vector \mathbf{v} . Note that the (k, h) th block of matrix Q , the (h, h') th block of matrix $Q^\top Q$, and the h th block of vector $Q^\top Q\mathbf{v}$ are given by $\nu_{kh} \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$,

3. Statistical test on the number of biclusters in a latent block model

$\omega_{hh'} \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$, and $\zeta_h \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$, respectively. Let $\mathbf{u}^{(h)} \in \mathbb{R}^p$ be a vector whose entries in the h th column cluster are $\frac{1}{\sqrt{p_h}}$ and whose all the other entries are zero.

Here, Theorem 2.1.2 (i.e., the *delocalization property*) holds for the eigenvectors $\{\mathbf{v}_j\}$ of matrix $Z^\top Z$ in our case, where we assume that $n, p \propto m$ and that each entry of Z is independently generated from a distribution with zero mean and unit variance that satisfies the sub-exponential condition. Therefore, we have $|\mathbf{v}^\top \mathbf{u}^{(h)}| = O_p(m^{-\frac{1}{2}+\epsilon})$ for all $\epsilon > 0$. Since $Q^\top Q \mathbf{v} = \sum_{h=1}^H \zeta_h \sqrt{p_h} \mathbf{u}^{(h)}$ and $\nu_{kh} = O_p(\frac{1}{m})$, $\omega_{hh'} = O_p(\frac{1}{m})$, $\zeta_h = \sum_{h'=1}^H \omega_{hh'} \sqrt{p_{h'}} \mathbf{v}^\top \mathbf{u}^{(h')} = O_p(m^{-1+\epsilon})$ for all $\epsilon > 0$ by definition, the following equation holds:

$$\begin{aligned} \|Q \mathbf{v}\|^2 &= \sum_{h=1}^H \zeta_h \sqrt{p_h} \mathbf{v}^\top \mathbf{u}^{(h)} = O_p(m^{-1+2\epsilon}) \\ \iff \|Q \mathbf{v}\| &= O_p(m^{-\frac{1}{2}+\epsilon}), \quad \text{for all } \epsilon > 0. \end{aligned} \quad (3.69)$$

Similarly, (h, h) th block of matrix $(\underline{Q}^{(k,h)})^\top \underline{Q}^{(k,h)}$ is $n_k \nu_{kh}^2 \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$, and its all the other entries are zero, which results in that $(\underline{Q}^{(k,h)})^\top \underline{Q}^{(k,h)} \mathbf{v} = n_k \nu_{kh}^2 \left(\sum_{j \in J_h} v_j \right) \sqrt{p_h} \mathbf{u}^{(h)}$. Therefore, we have

$$\begin{aligned} \|\underline{Q}^{(k,h)} \mathbf{v}\|^2 &= n_k \nu_{kh}^2 \left(\sum_{j \in J_h} v_j \right) \sqrt{p_h} \mathbf{v}^\top \mathbf{u}^{(h)} = n_k \nu_{kh}^2 \left(\sqrt{p_h} \mathbf{v}^\top \mathbf{u}^{(h)} \right)^2 \\ &= O_p(m^{-1+2\epsilon}) \\ \iff \|\underline{Q}^{(k,h)} \mathbf{v}\| &= O_p(m^{-\frac{1}{2}+\epsilon}), \quad \text{for all } \epsilon > 0. \end{aligned} \quad (3.70)$$

Moreover, from Lemma A1, we have

$$|1 - \tau_{kh}| = O_p\left(\frac{1}{m}\right). \quad (3.71)$$

By substituting (3.69), (3.70), and (3.71) into (3.67) and by setting $\epsilon < \frac{1}{2}$, we have

$$\begin{aligned} \tilde{\lambda}_1 &\geq \lambda_1 - O_p(m^\epsilon) - O_p\left(m^{\frac{1}{2}}\right) \left[\sum_{k=1}^K \sum_{h=1}^H O_p(m^{-1}) O_p\left(m^{\frac{1}{2}}\right) \right] \\ &= \lambda_1 - O_p(m^\epsilon), \end{aligned} \quad (3.72)$$

3. Statistical test on the number of biclusters in a latent block model

which concludes the proof. Here, we used the assumption that K and H are fixed constants. \square

Lemma B2. *Let λ_1 and $\tilde{\lambda}_1$, respectively, be the maximum eigenvalues of matrices $Z^\top Z$ and $\tilde{Z}^\top \tilde{Z}$. Then, the following equation holds:*

$$\tilde{\lambda}_1 \leq \lambda_1 + O_p\left(m^{\frac{2}{7}+\epsilon}\right), \quad \text{for all } \epsilon \in \left(0, \frac{2}{7}\right). \quad (3.73)$$

Proof. We use the same notations as in Lemmas B1. Let $\{\lambda_j\}$ and $\{\mathbf{v}_j\}$, respectively, be the sets of the eigenvalues and the corresponding normalized eigenvectors (i.e., $\|\mathbf{v}_j\| = 1$ for all j) of matrix $Z^\top Z$, where $j = 1, \dots, p$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. Let $\tilde{\mathbf{v}}^{(h)} \in \mathbb{R}^{p_h}$ be a subvector of $\tilde{\mathbf{v}}$ in the h th column cluster.

Since $Z^\top Z$ is a symmetric matrix, its eigenvectors $\{\mathbf{v}_j\}$ form an orthonormal system, and thus there exists a unique set of coefficients $\{c_j\}$ and t ($1 \leq t \leq p$) that satisfies

$$\tilde{\mathbf{v}} = \sum_{j=1}^p c_j \mathbf{v}_j = \tilde{\mathbf{v}}_1 + \tilde{\mathbf{v}}_2, \quad (3.74)$$

where

$$\begin{aligned} \tilde{\mathbf{v}}_1 &\equiv \sum_{j=1}^t c_j \mathbf{v}_j, & \tilde{\mathbf{v}}_2 &\equiv \sum_{j=t+1}^p c_j \mathbf{v}_j, \\ \lambda_t &\geq \lambda_1 - n^d, & \lambda_{t+1} &< \lambda_1 - n^d, & d &= \frac{5}{7}. \end{aligned} \quad (3.75)$$

Therefore, by using (3.71), the following equation holds:

$$\begin{aligned} \tilde{\lambda}_1 &= \tilde{\mathbf{v}}^\top \tilde{Z}^\top \tilde{Z} \tilde{\mathbf{v}} = \left\| \sum_{k=1}^K \sum_{h=1}^H \left[\underline{Z}^{(k,h)} + (\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \right\|^2 \\ &= \left\| \left\{ Z + \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \right\} \tilde{\mathbf{v}} \right\|^2 \\ &= \|Z \tilde{\mathbf{v}}\|^2 + 2 \tilde{\mathbf{v}}^\top Z^\top \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \\ &\quad + \left\| \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \right\|^2 \\ &\leq \|Z \tilde{\mathbf{v}}\|^2 + 2\sqrt{\lambda_1} \sum_{k=1}^K \sum_{h=1}^H |\tau_{kh} - 1| \|\underline{Z}^{(k,h)} \tilde{\mathbf{v}}\| - 2 \tilde{\mathbf{v}}^\top Z^\top \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
& + \left\| \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \right\|^2 \\
& = \|\underline{Z} \tilde{\mathbf{v}}\|^2 + 2\sqrt{\lambda_1} \sum_{k=1}^K \sum_{h=1}^H |\tau_{kh} - 1| \|\underline{Z}^{(k,h)} \tilde{\mathbf{v}}^{(h)}\| - 2\tilde{\mathbf{v}}^\top \underline{Z}^\top \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
& + \left\| \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \right\|^2 \\
& = \|\underline{Z} \tilde{\mathbf{v}}\|^2 + O_p(\sqrt{m}) \sum_{k=1}^K \sum_{h=1}^H O_p\left(\frac{1}{m}\right) O_p(\sqrt{m}) - 2\tilde{\mathbf{v}}^\top \underline{Z}^\top \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
& + \left\| \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \right\|^2 \\
& = \|\underline{Z} \tilde{\mathbf{v}}\|^2 + O_p(1) - 2\tilde{\mathbf{v}}^\top \underline{Z}^\top \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
& + \left\| \sum_{k=1}^K \sum_{h=1}^H \left[(\tau_{kh} - 1) \underline{Z}^{(k,h)} - \tau_{kh} \underline{Q}^{(k,h)} \right] \tilde{\mathbf{v}} \right\|^2 \quad (\because K \text{ and } H \text{ are fixed constants}) \\
& \leq \|\underline{Z} \tilde{\mathbf{v}}\|^2 + O_p(1) - 2\tilde{\mathbf{v}}^\top \underline{Z}^\top \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
& + \left(\left\| \sum_{k=1}^K \sum_{h=1}^H (\tau_{kh} - 1) \underline{Z}^{(k,h)} \tilde{\mathbf{v}} \right\| + \left\| \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \right\| \right)^2 \\
& \leq \|\underline{Z} \tilde{\mathbf{v}}\|^2 + O_p(1) - 2\tilde{\mathbf{v}}^\top \underline{Z}^\top \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
& + \left(\sum_{k=1}^K \sum_{h=1}^H |\tau_{kh} - 1| \|\underline{Z}^{(k,h)} \tilde{\mathbf{v}}\| + \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \|\underline{Q}^{(k,h)} \tilde{\mathbf{v}}\| \right)^2 \\
& = \|\underline{Z} \tilde{\mathbf{v}}\|^2 + O_p(1) - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top \underline{Z}^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
& + \left[O_p\left(\frac{1}{\sqrt{m}}\right) + \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \|\underline{Q}^{(k,h)} \tilde{\mathbf{v}}\| \right]^2 \quad (\because K \text{ and } H \text{ are fixed constants}). \quad (3.76)
\end{aligned}$$

As for the last term in (3.76), the following equation holds:

$$\|\underline{Q}^{(k,h)} \tilde{\mathbf{v}}\|^2 = \tilde{\mathbf{v}}^\top (\underline{Q}^{(k,h)})^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} = \sum_{j=1}^p \sum_{j'=1}^p c_j c_{j'} \mathbf{v}_j^\top (\underline{Q}^{(k,h)})^\top \underline{Q}^{(k,h)} \mathbf{v}_{j'}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&= O_p\left(\frac{1}{m}\right) \sum_{j=1}^p \sum_{j'=1}^p c_j c_{j'} \mathbf{v}_j^\top p_h(\mathbf{v}_{j'}^\top \mathbf{u}^{(h)}) \mathbf{u}^{(h)} \\
&= O_p(1) \left[\sum_{j=1}^p c_j (\mathbf{v}_j^\top \mathbf{u}^{(h)}) \right]^2 \leq O_p(1) \left[\sqrt{\sum_{j=1}^p c_j^2} \sqrt{\sum_{j=1}^p (\mathbf{v}_j^\top \mathbf{u}^{(h)})^2} \right]^2 \\
&= O_p(1) \left(\sum_{j=1}^p c_j^2 \right) \left[\sum_{j=1}^p (\mathbf{v}_j^\top \mathbf{u}^{(h)})^2 \right] = O_p(1) \|\tilde{\mathbf{v}}\|^2 \left[\sum_{j=1}^p (\mathbf{v}_j^\top \mathbf{u}^{(h)})^2 \right] \\
&= O_p(1) \left[\sum_{j=1}^p (\mathbf{v}_j^\top \mathbf{u}^{(h)})^2 \right] \\
&= O_p(m^{2\epsilon}) \quad (\because \text{Theorem 2.1.2}) \tag{3.77}
\end{aligned}$$

$$\iff \|\underline{Q}^{(k,h)} \tilde{\mathbf{v}}\| = O_p(m^\epsilon), \quad \text{for all } \epsilon > 0, \tag{3.78}$$

where $\mathbf{u}^{(h)} \in \mathbb{R}^p$ is a vector whose elements in the h th column cluster is $\frac{1}{\sqrt{p_h}}$ and whose all the other elements are zero. In the last equation in (3.77), we used the delocalization property of $\{\mathbf{v}_j\}$, which are eigenvectors of matrix $Z^\top Z$. By substituting (3.78) into (3.76) and using (3.71) and the assumption that K and H are fixed constants, we have

$$\tilde{\lambda}_1 \leq \|Z\tilde{\mathbf{v}}\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} + O_p(m^{2\epsilon}), \quad \text{for all } \epsilon > 0. \tag{3.79}$$

Here, by definition in (3.75), the following equation holds:

$$\begin{aligned}
&\|Z\tilde{\mathbf{v}}\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
&= (\tilde{\mathbf{v}}_1 + \tilde{\mathbf{v}}_2)^\top Z^\top Z (\tilde{\mathbf{v}}_1 + \tilde{\mathbf{v}}_2) - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
&= \tilde{\mathbf{v}}_1^\top Z^\top Z \tilde{\mathbf{v}}_1 + \tilde{\mathbf{v}}_2^\top Z^\top Z \tilde{\mathbf{v}}_2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
&= \left(\sum_{j=1}^t c_j \mathbf{v}_j \right)^\top \sum_{j=1}^t c_j \lambda_j \mathbf{v}_j + \left(\sum_{j=t+1}^p c_j \mathbf{v}_j \right)^\top \sum_{j=t+1}^p c_j \lambda_j \mathbf{v}_j - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \\
&= \sum_{j=1}^t c_j^2 \lambda_j + \sum_{j=t+1}^p c_j^2 \lambda_j - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} \quad (\because \|\mathbf{v}_j\| = 1 \text{ for all } j) \\
&\leq \lambda_1 \sum_{j=1}^t c_j^2 + (\lambda_1 - n^d) \sum_{j=t+1}^p c_j^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}
\end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&= \lambda_1 \|\tilde{\mathbf{v}}\|^2 - n^d \|\tilde{\mathbf{v}}_2\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} (\tilde{\mathbf{v}}_1 + \tilde{\mathbf{v}}_2) \\
&= \lambda_1 - n^d \|\tilde{\mathbf{v}}_2\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_1 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_2. \quad (3.80)
\end{aligned}$$

The third term in (3.80) can be upper bounded as follows:

$$\begin{aligned}
&- \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_1 \leq |\tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_1| = \left| \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \left(\sum_{j=1}^t c_j \mathbf{v}_j \right) \right| \\
&= \left| \sum_{j=1}^t c_j \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \mathbf{v}_j \right| \\
&= O_p \left(\frac{1}{m} \right) \left| \sum_{j=1}^t c_j \left[\sqrt{n_k} (\tilde{\mathbf{w}}^{(k)})^\top Z \tilde{\mathbf{v}} \right]^\top \sqrt{p_h} (\tilde{\mathbf{u}}^{(h)})^\top \mathbf{v}_j \right| \\
&= O_p(1) \left| (\tilde{\mathbf{v}}^\top Z^\top \tilde{\mathbf{w}}^{(k)}) \sum_{j=1}^t c_j (\tilde{\mathbf{u}}^{(h)})^\top \mathbf{v}_j \right| \\
&\leq O_p(1) \|Z \tilde{\mathbf{v}}\| \|\tilde{\mathbf{w}}^{(k)}\| \sum_{j=1}^t c_j (\tilde{\mathbf{u}}^{(h)})^\top \mathbf{v}_j \leq O_p(1) \|Z\|_{\text{op}} \sum_{j=1}^t c_j (\tilde{\mathbf{u}}^{(h)})^\top \mathbf{v}_j \\
&= O_p(\sqrt{m}) \left| \sum_{j=1}^t c_j (\tilde{\mathbf{u}}^{(h)})^\top \mathbf{v}_j \right| \leq O_p(\sqrt{m}) \sqrt{\sum_{j=1}^t c_j^2} \sqrt{\sum_{j=1}^t \left[(\tilde{\mathbf{u}}^{(h)})^\top \mathbf{v}_j \right]^2} \\
&= O_p(\sqrt{m}) \|\tilde{\mathbf{v}}\| \sqrt{t} O_p \left(m^{-\frac{1}{2} + \epsilon} \right) \quad (\because \text{Theorem 2.1.2}) \\
&= \sqrt{t} O_p(m^\epsilon), \quad \text{for all } \epsilon > 0, \quad (3.81)
\end{aligned}$$

where $\tilde{\mathbf{w}}^{(k)} \in \mathbb{R}^n$ is a vector whose elements in the k th row cluster is $\frac{1}{\sqrt{n_k}}$ and whose all the other elements are zero. Here we used the delocalization property of $\{\mathbf{v}_j\}$, which are eigenvectors of matrix $Z^\top Z$.

The fourth term in (3.80) can also be upper bounded as follows:

$$\begin{aligned}
&- \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_2 \leq |\tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_2| \leq \|\tilde{\mathbf{v}}\| \|Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_2\| = \|Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}_2\| \\
&\leq \|Z^\top \underline{Q}^{(k,h)}\|_{\text{op}} \|\tilde{\mathbf{v}}_2\| \leq \|Z\|_{\text{op}} \|\underline{Q}^{(k,h)}\|_{\text{F}} \|\tilde{\mathbf{v}}_2\| \\
&= O_p(\sqrt{m}) O_p \left(\frac{1}{m} \right) \sqrt{n_k p_h} \|\tilde{\mathbf{v}}_2\| = \|\tilde{\mathbf{v}}_2\| O_p(\sqrt{m}). \quad (3.82)
\end{aligned}$$

By substituting (3.81) and (3.82) into (3.80), we have

$$\|Z \tilde{\mathbf{v}}\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&\leq \lambda_1 - n^d \|\tilde{\mathbf{v}}_2\|^2 + 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \sqrt{t} O_p(m^\epsilon) + 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \|\tilde{\mathbf{v}}_2\| O_p(\sqrt{m}) \\
&= \lambda_1 - n^d \|\tilde{\mathbf{v}}_2\|^2 + \sqrt{t} O_p(m^\epsilon) + \|\tilde{\mathbf{v}}_2\| O_p(\sqrt{m}). \tag{3.83}
\end{aligned}$$

In the last equation of (3.83), we used (3.71) and the assumption that K and H are fixed constants.

Let $\nu_j \equiv \frac{1}{n} \lambda_j$ be a normalized eigenvalue of matrix $Z^\top Z$. Note that t in (3.75) is the number of normalized eigenvalues $\{\nu_j\}$ that satisfy $\nu_j \geq \nu_1 - n^{d-1}$. We also define the following variables:

$$\nu_+ \equiv \left(1 + \sqrt{\frac{p}{n}}\right)^2, \quad \nu_- \equiv \left(1 - \sqrt{\frac{p}{n}}\right)^2, \quad \epsilon_1 \equiv \nu_+ - \nu_1. \tag{3.84}$$

From (4.1) of [122], $|\epsilon_1| = O_p(\phi^C m^{-\frac{2}{3}})$ holds for some constant $C > 0$, where $\phi \equiv (\log p)^{\log \log p}$. Since $\phi = o(m^{\tilde{\epsilon}_0})$ holds for any $\tilde{\epsilon}_0 > 0$, we have $|\epsilon_1| = O_p(m^{-\frac{2}{3} + \epsilon_0})$ for any $\epsilon_0 > 0$.

From (3.7) of [122],

$$\left| \bar{n} - \frac{t}{p} \right| = O_p(m^{-1 + \epsilon_2}), \quad \forall \epsilon_2 > 0, \tag{3.85}$$

where

$$\begin{aligned}
\bar{n} &= \int_{\nu_1 - n^{d-1}}^{\infty} q(x) dx, \\
q(x) &= \frac{1}{2\pi} \frac{n}{p} \frac{\sqrt{\max\{(\nu_+ - x)(x - \nu_-), 0\}}}{x}. \tag{3.86}
\end{aligned}$$

By setting $\epsilon_0 < d - \frac{1}{3}$, we have

$$\begin{aligned}
q(\nu_1 - n^{d-1}) &= q(\nu_+ - n^{d-1} - \epsilon_1) \\
&= \frac{\sqrt{\nu_+ - \nu_-}}{\nu_+} \left[n^{\frac{d-1}{2}} + O_p\left(m^{-\frac{1}{3} + \frac{\epsilon_0}{2}}\right) \right] \left[1 + O\left(m^{\frac{d-1}{2}}\right) + O_p\left(m^{-\frac{1}{3} + \frac{\epsilon_0}{2}}\right) \right] \\
&= \frac{\sqrt{\nu_+ - \nu_-}}{\nu_+} n^{\frac{d-1}{2}} + O_p\left(m^{\max\{\frac{d-1}{2}, -\frac{1}{3} + \frac{\epsilon_0}{2}\}}\right) \\
&= \frac{\sqrt{\nu_+ - \nu_-}}{\nu_+} n^{\frac{d-1}{2}} + O_p\left(m^{\frac{d-1}{2}}\right) \quad \left(\because \epsilon_0 < d - \frac{1}{3}\right). \tag{3.87}
\end{aligned}$$

From (3.87) and the fact that $|\epsilon_1| = O_p(m^{-\frac{2}{3} + \epsilon_0})$ for any $\epsilon_0 > 0$, by setting $\epsilon_0 < d - \frac{1}{3}$, the following equation holds:

$$\bar{n} = \int_{\nu_1 - n^{d-1}}^{\infty} q(x) dx \leq \left| \int_{\nu_1 - n^{d-1}}^{\nu_+} q(x) dx \right| + \left| \int_{\nu_+}^{\infty} q(x) dx \right|$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&= \left| \int_{\nu_1 - n^{d-1}}^{\nu_+} q(x) dx \right| \leq |\epsilon_1 + n^{d-1}| q(\nu_1 - n^{d-1}) \\
&= O_p(m^{d-1}) O_p\left(m^{\frac{d-1}{2}}\right) = O_p\left(m^{\frac{3(d-1)}{2}}\right). \tag{3.88}
\end{aligned}$$

Therefore, from (3.85), by setting $\epsilon_2 < \frac{3}{2}d - \frac{1}{2}$, we have

$$t = O_p\left(m^{\frac{3}{2}d - \frac{1}{2}}\right). \tag{3.89}$$

By combining (3.83) and the assumption in (3.75) that $d = \frac{5}{7}$, the following equation holds for all $\epsilon > 0$:

$$\begin{aligned}
\|Z\tilde{\mathbf{v}}\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} &\leq \lambda_1 + \|\tilde{\mathbf{v}}_2\| \left[n^{\frac{1}{2}}\varpi - n^d \|\tilde{\mathbf{v}}_2\| \right] + O_p\left(m^{\frac{2}{7}+\epsilon}\right), \\
\varpi \equiv n^{-\frac{1}{2}} \|Z\|_{\text{op}} \|\underline{Q}^{(k,h)}\|_{\text{F}} &= O_p(1) \quad (\because (3.82)). \tag{3.90}
\end{aligned}$$

Here, we consider the following two patterns: (a) If $n^{\frac{1}{2}}\varpi - n^d \|\tilde{\mathbf{v}}_2\| \leq 0$, from (3.90), we have

$$\|Z\tilde{\mathbf{v}}\|^2 - 2 \sum_{k=1}^K \sum_{h=1}^H \tau_{kh} \tilde{\mathbf{v}}^\top Z^\top \underline{Q}^{(k,h)} \tilde{\mathbf{v}} = \lambda_1 + O_p\left(m^{\frac{2}{7}+\epsilon}\right). \tag{3.91}$$

(b) If $n^{\frac{1}{2}}\varpi - n^d \|\tilde{\mathbf{v}}_2\| > 0$, we have $\|\tilde{\mathbf{v}}_2\| < n^{\frac{1}{2}-d}\varpi$ and thus

$$\|\tilde{\mathbf{v}}_2\| \left[n^{\frac{1}{2}}\varpi - n^d \|\tilde{\mathbf{v}}_2\| \right] \leq n^{1-d}\varpi^2. \tag{3.92}$$

By assumption in (3.75) that $d = \frac{5}{7}$, we have $n^{1-d}\varpi^2 = O_p\left(m^{\frac{2}{7}}\right)$ and thus (3.91) holds.

In summary, (3.91) always holds. By combining this fact and (3.79), we have

$$\tilde{\lambda}_1 \leq \lambda_1 + O_p\left(m^{\frac{2}{7}+\epsilon}\right) + O_p\left(m^{2\epsilon}\right), \quad \text{for all } \epsilon > 0. \tag{3.93}$$

By setting $\epsilon < \frac{2}{7}$, we finally obtain

$$\tilde{\lambda}_1 \leq \lambda_1 + O_p\left(m^{\frac{2}{7}+\epsilon}\right), \quad \text{for all } \epsilon \in \left(0, \frac{2}{7}\right), \tag{3.94}$$

which concludes the proof. \square

Lemma B3. Let λ_1 and $\hat{\lambda}_1$, respectively, be the maximum eigenvalues of matrices $Z^\top Z$ and $\hat{Z}^\top \hat{Z}$ (i.e., $\|Z\|_{\text{op}}^2$ and $\|\hat{Z}\|_{\text{op}}^2$, respectively). Then, for all $\epsilon \in (0, \frac{2}{7})$,

$$\frac{|\lambda_1 - \hat{\lambda}_1|}{b^{\text{TW}}} = O_p\left(m^{-\frac{1}{21}+\epsilon}\right), \tag{3.95}$$

where b^{TW} is defined as in (2.3).

3. Statistical test on the number of biclusters in a latent block model

Proof. From Lemma B1 and B2, we have already shown that the following equation holds for all $\epsilon \in (0, \frac{2}{7})$:

$$|\lambda_1 - \tilde{\lambda}_1| = O_p \left(m^{\frac{2}{7} + \epsilon} \right) \iff \frac{|\lambda_1 - \tilde{\lambda}_1|}{b^{\text{TW}}} = O_p \left(m^{-\frac{1}{21} + \epsilon} \right). \quad (3.96)$$

We consider the joint probability of the event \mathcal{F}_m that $\tilde{Z} = \hat{Z}$ holds and the event $\mathcal{G}_{m,C}$ that $\frac{|\lambda_1 - \tilde{\lambda}_1|}{b^{\text{TW}}} \leq C m^{-\frac{1}{21} + \epsilon}$ holds. Such a joint probability satisfies the following inequality:

$$\Pr(\mathcal{F}_m \cap \mathcal{G}_{m,C}) \geq 1 - \Pr(\mathcal{F}_m^C) - \Pr(\mathcal{G}_{m,C}^C), \quad (3.97)$$

where \mathcal{A}^C is the complement of event \mathcal{A} . The consistency assumption (vi) guarantees that if $(K_0, H_0) = (K, H)$, $\Pr(\mathcal{F}_m^C)$ converges to 0 in the limit of $m \rightarrow \infty$. By combining this fact with (3.96), for all $\tilde{\epsilon} > 0$, there exist $C > 0$ and $M > 0$ such that for all $m \geq M$, $\Pr(\mathcal{F}_m \cap \mathcal{G}_{m,C}) \geq 1 - \tilde{\epsilon}$ holds, which results in (3.95). \square

3.C Proof of $\hat{\sigma}^* = O_p(1)$ in unrealizable case

Proof. Throughout the proof, we use the following notations:

- $A^{(k,h)}$, $P^{(k,h)}$, and $Z^{(k,h)}$, respectively, are the (k, h) th **null** blocks of matrices A , P , and Z .
- $\underline{A}^{(k,h)}$, $\underline{P}^{(k,h)}$, and $\hat{P}^{(k,h)}$, respectively, are the (k, h) th **estimated** blocks of matrices A , P , and \hat{P} .
- We denote the row and column sizes of the (k, h) th **estimated** block as \underline{n}_k and \underline{p}_h , respectively.
- (k_1, h_1) is the set of row and column cluster indices of submatrix \bar{X} in the **estimated** block structure.

As for the order of the estimated standard deviation $\hat{\sigma}^*$, we have $\hat{\sigma}^* = \hat{S}_{k_1 h_1}$. Note that the block size (\bar{n}_1, \bar{p}_1) of submatrix \bar{X} is at least $(n_{\min}/K_0) \times (p_{\min}/H_0)$. Therefore, we have

$$\begin{aligned} \hat{\sigma}^* &= \hat{S}_{k_1 h_1} = \frac{1}{\sqrt{\underline{n}_{k_1} \underline{p}_{h_1}}} \|\underline{A}^{(k_1, h_1)} - \hat{P}^{(k_1, h_1)}\|_{\text{F}} \\ &\leq \frac{1}{\sqrt{\bar{n}_1 \bar{p}_1}} \|\underline{A}^{(k_1, h_1)} - \hat{P}^{(k_1, h_1)}\|_{\text{F}} \leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \|\underline{A}^{(k_1, h_1)} - \hat{P}^{(k_1, h_1)}\|_{\text{F}} \\ &\leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \|A - \hat{P}\|_{\text{F}} = \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \|A - P + P - \hat{P}\|_{\text{F}} \end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&\leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \left(\|A - P\|_{\text{F}} + \|P - \hat{P}\|_{\text{F}} \right) \\
&= \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \left(\sqrt{\sum_{k=1}^K \sum_{h=1}^H \|A^{(k,h)} - P^{(k,h)}\|_{\text{F}}^2} + \|P - \hat{P}\|_{\text{F}} \right) \\
&= \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \left(\sqrt{\sum_{k=1}^K \sum_{h=1}^H S_{kh}^2 \|Z^{(k,h)}\|_{\text{F}}^2} + \|P - \hat{P}\|_{\text{F}} \right) \\
&\leq \sqrt{\frac{K_0 H_0}{n_{\min} p_{\min}}} \left[\sqrt{KH} \left(\max_{k=1, \dots, K, h=1, \dots, H} S_{kh} \right) \|Z\|_{\text{F}} + \|P - \hat{P}\|_{\text{F}} \right]. \quad (3.98)
\end{aligned}$$

Here, for all (i, j) , $\left(Z_{ij}^{(k,h)}\right)^2$ independently follows the same distribution, and $\mathbb{E} \left[\left(Z_{ij}^{(k,h)}\right)^2 \right] = 1$. We also have $\mathbb{V} \left[\left(Z_{ij}^{(k,h)}\right)^2 \right] = \mathbb{E} \left[\left(Z_{ij}^{(k,h)}\right)^4 \right] - 1 < \infty$, since we have assumed that $\mathbb{E} \left[\left(Z_{ij}^{(k,h)}\right)^4 \right] < \infty$ from the sub-exponential assumption. Therefore, from the central limit theorem and Prokhorov's theorem [146], we have $\frac{1}{\sqrt{n_k p_h}} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left[\left(Z_{ij}^{(k,h)}\right)^2 - 1 \right] = O_p(1)$. In other words, the following equation holds: $\sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(Z_{ij}^{(k,h)}\right)^2 = n_k p_h + O_p(m) = O_p(m^2)$. Based on this result, we obtain

$$\begin{aligned}
\|Z\|_{\text{F}} &= \sqrt{\sum_{k=1}^K \sum_{h=1}^H \|Z^{(k,h)}\|_{\text{F}}^2} = \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} \left(Z_{ij}^{(k,h)}\right)^2} \\
&= \sqrt{\sum_{k=1}^K \sum_{h=1}^H O_p(m^2)} = O_p(m). \quad (3.99)
\end{aligned}$$

Here, we used the assumption that K and H are fixed constants.

Furthermore, we have

$$\begin{aligned}
\|P - \hat{P}\|_{\text{F}} &= \sqrt{\sum_{i=1}^n \sum_{j=1}^p \left(P_{ij} - \hat{P}_{ij}\right)^2} = \sqrt{\sum_{i=1}^n \sum_{j=1}^p \left(P_{ij} - \bar{P}_{ij} + \bar{P}_{ij} - \hat{P}_{ij}\right)^2} \\
&\leq \sqrt{\sum_{i=1}^n \sum_{j=1}^p \left(|P_{ij} - \bar{P}_{ij}| + |\bar{P}_{ij} - \hat{P}_{ij}|\right)^2}
\end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&\leq \sqrt{\sum_{i=1}^n \sum_{j=1}^p \left[\left(\max_{i'=1, \dots, n, j'=1, \dots, p} |P_{i'j'} - \bar{P}_{i'j'}| \right) + \left| \bar{P}_{ij} - \hat{P}_{ij} \right| \right]^2} \\
&\leq \sqrt{\sum_{i=1}^n \sum_{j=1}^p \left[\left(\max_{\substack{k=1, \dots, K, h=1, \dots, H, \\ k'=1, \dots, K, h'=1, \dots, H}} |B_{kh} - B_{k'h'}| \right) + \left| \bar{P}_{ij} - \hat{P}_{ij} \right| \right]^2} \\
&\leq \sqrt{\sum_{i=1}^n \sum_{j=1}^p \left(\max_{\substack{k=1, \dots, K, h=1, \dots, H, \\ k'=1, \dots, K, h'=1, \dots, H}} |B_{kh} - B_{k'h'}| + \max_{k=1, \dots, K_0, h=1, \dots, H_0} |\bar{B}_{kh} - \hat{B}_{kh}| \right)^2} \\
&= \sqrt{np} \left(\max_{\substack{k=1, \dots, K, h=1, \dots, H, \\ k'=1, \dots, K, h'=1, \dots, H}} |B_{kh} - B_{k'h'}| + \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} |\bar{B}_{kh} - \hat{B}_{kh}| \right) \\
&\leq \sqrt{np} \left[\max_{\substack{k=1, \dots, K, h=1, \dots, H, \\ k'=1, \dots, K, h'=1, \dots, H}} |B_{kh} - B_{k'h'}| + O_p \left(\frac{1}{\sqrt{m}} \right) \right] \\
&= O_p(m). \tag{3.100}
\end{aligned}$$

Here, to derive the last inequality in (3.100), we used (3.29) and the assumption that K_0 and H_0 are fixed constants. In the final equation, we used the fact that $\max_{\substack{k=1, \dots, K, h=1, \dots, H, \\ k'=1, \dots, K, h'=1, \dots, H}} |B_{kh} - B_{k'h'}|$ is bounded by a finite constant.

By combining (3.98), (3.99), and (3.100), we obtain $\hat{\sigma}^* = O_p(1)$. \square

3.D Proof of the asymptotic ICL in the Bernoulli case

Proof. From Lemma 4.2 in [76], the resulting asymptotic ICL is given by

$$\begin{aligned}
\text{ICL}(K_0, H_0) &= \max_{\pi, \rho, B} \log p(A, \hat{g}^{(1)}, \hat{g}^{(2)} | \pi, \rho, B) \\
&\quad - \frac{K_0 - 1}{2} \log n - \frac{H_0 - 1}{2} \log p - \frac{K_0 H_0}{2} \log(np). \tag{3.101}
\end{aligned}$$

In regard to the first term in (3.101), we consider the following optimization problem:

$$\begin{aligned}
&\max_{\pi, \rho, B} \log p(A, \hat{g}^{(1)}, \hat{g}^{(2)} | \pi, \rho, B), \\
\text{s.t. } &\sum_{k=1}^{K_0} \pi_k = 1, \pi_k \geq 0 \text{ for all } k, \sum_{h=1}^{H_0} \rho_h = 1, \rho_h \geq 0 \text{ for all } h, \\
&0 \leq B_{kh} \leq 1 \text{ for all } (k, h). \tag{3.102}
\end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

The above problem is solved with the Lagrangian undetermined multiplier method, which employs

$$\begin{aligned}
f &\equiv \log p(A, \hat{g}^{(1)}, \hat{g}^{(2)} | \pi, \rho, B) - \xi_1 \sum_{k=1}^{K_0} \pi_k - \xi_2 \sum_{h=1}^{H_0} \rho_h, \\
&= \sum_{k=1}^{K_0} |I_k| \log \pi_k + \sum_{h=1}^{H_0} |J_h| \log \rho_h - \xi_1 \sum_{k=1}^{K_0} \pi_k - \xi_2 \sum_{h=1}^{H_0} \rho_h \\
&+ \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \sum_{i \in I_k, j \in J_h} [A_{ij} \log B_{kh} + (1 - A_{ij}) \log (1 - B_{kh})], \tag{3.103}
\end{aligned}$$

$$\frac{\partial f}{\partial \pi_k} = \frac{\partial f}{\partial \rho_h} = \frac{\partial f}{\partial B_{kh}} = 0 \text{ for all } k, h. \tag{3.104}$$

By substituting (3.103) into (3.104), we have

$$\begin{aligned}
\frac{|I_k|}{\pi_k} = \xi_1, \quad \frac{|J_h|}{\rho_h} = \xi_2, \quad \sum_{i \in I_k, j \in J_h} \left[\frac{A_{ij}}{B_{kh}} - \frac{1 - A_{ij}}{1 - B_{kh}} \right] &= 0 \\
\iff \pi_k = \frac{|I_k|}{\xi_1}, \quad \rho_h = \frac{|J_h|}{\xi_2}, \quad B_{kh} = \frac{\sum_{i \in I_k, j \in J_h} A_{ij}}{|I_k| |J_h|}, &\tag{3.105}
\end{aligned}$$

for all (k, h) . In regard to $\{\pi_k\}$ and $\{\rho_h\}$, from the conditions in (3.102), $\sum_{k=1}^{K_0} |I_k| = \xi_1$ and $\sum_{h=1}^{H_0} |J_h| = \xi_2$ hold and thus we finally have

$$\pi_k = \frac{|I_k|}{\sum_{k=1}^{K_0} |I_k|}, \quad \rho_h = \frac{|J_h|}{\sum_{h=1}^{H_0} |J_h|}. \tag{3.106}$$

We can easily check that the solutions of (3.105) and (3.106) satisfy all the conditions in (3.102).

Finally, by substituting the above results into (3.101), we have

$$\begin{aligned}
\text{ICL}(K_0, H_0) &= \sum_{k=1}^{K_0} |I_k| \log \left(\frac{|I_k|}{\sum_{k=1}^{K_0} |I_k|} \right) + \sum_{h=1}^{H_0} |J_h| \log \left(\frac{|J_h|}{\sum_{h=1}^{H_0} |J_h|} \right) \\
&+ \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \left(\sum_{i \in I_k, j \in J_h} A_{ij} \right) \log \left(\frac{\sum_{i \in I_k, j \in J_h} A_{ij}}{|I_k| |J_h|} \right) \\
&+ \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \left(|I_k| |J_h| - \sum_{i \in I_k, j \in J_h} A_{ij} \right) \log \left(1 - \frac{\sum_{i \in I_k, j \in J_h} A_{ij}}{|I_k| |J_h|} \right) \\
&- \frac{K_0 - 1}{2} \log n - \frac{H_0 - 1}{2} \log p - \frac{K_0 H_0}{2} \log(np)
\end{aligned}$$

3. Statistical test on the number of biclusters in a latent block model

$$\begin{aligned}
&= \sum_{k=1}^{K_0} |I_k| \log \left(\frac{|I_k|}{n} \right) + \sum_{h=1}^{H_0} |J_h| \log \left(\frac{|J_h|}{p} \right) \\
&+ \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} |I_k| |J_h| \left[\hat{B}_{kh} \log \hat{B}_{kh} + (1 - \hat{B}_{kh}) \log (1 - \hat{B}_{kh}) \right] \\
&- \frac{K_0 - 1}{2} \log n - \frac{H_0 - 1}{2} \log p - \frac{K_0 H_0}{2} \log(np). \tag{3.107}
\end{aligned}$$

Note that we have defined \hat{B}_{kh} as in (3.9). \square

3.E Jarque–Bera test for selecting the cluster numbers

For Gaussian LBMs, we can use Jarque–Bera test [12, 22, 71] instead of the proposed one for determining the row and column cluster numbers. We define the sample skewness $\tilde{\theta}_{kh}$ and sample kurtosis $\tilde{\kappa}_{kh}$ of the (k, h) th **null** block as follows:

$$\begin{aligned}
\tilde{\theta} &= (\tilde{\theta}_{kh})_{1 \leq k \leq K, 1 \leq h \leq H}, & \tilde{\theta}_{kh} &\equiv \frac{\frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} (A_{ij} - \tilde{B}_{kh})^3}{\tilde{S}_{kh}^3}, \\
\tilde{\kappa} &= (\tilde{\kappa}_{kh})_{1 \leq k \leq K, 1 \leq h \leq H}, & \tilde{\kappa}_{kh} &\equiv \frac{\frac{1}{n_k p_h} \sum_{i=1}^{n_k} \sum_{j=1}^{p_h} (A_{ij} - \tilde{B}_{kh})^4}{\tilde{S}_{kh}^4}. \tag{3.108}
\end{aligned}$$

From the result of [22, 71], for each (k, h) th block, $\tilde{T}_{kh}^{\text{JB}} \equiv \frac{n_k p_h}{6} \left[\tilde{\theta}_{kh}^2 + \frac{1}{4} (\tilde{\kappa}_{kh} - 3)^2 \right]$ converges in law to the chi-squared distribution with 2 degrees of freedom. Since we assume that $\{A_{ij}\}$ are mutually independent, given the block structure, the statistic $\tilde{T}^{\text{JB}} \equiv \sum_{k=1}^K \sum_{h=1}^H \tilde{T}_{kh}^{\text{JB}}$ converges in law to the chi-squared distribution with $2KH$ degrees of freedom. Here, we used the additivity of the chi-squared distribution.

We also define the sample skewness $\hat{\theta}_{kh}$ and sample kurtosis $\hat{\kappa}_{kh}$ of the (k, h) th **estimated** block as follows:

$$\begin{aligned}
\hat{\theta} &= (\hat{\theta}_{kh})_{1 \leq k \leq K_0, 1 \leq h \leq H_0}, & \hat{\theta}_{kh} &\equiv \frac{\frac{1}{|I_k| |J_h|} \sum_{i=1}^{|I_k|} \sum_{j=1}^{|J_h|} (A_{ij} - \hat{B}_{kh})^3}{\hat{S}_{kh}^3}, \\
\hat{\kappa} &= (\hat{\kappa}_{kh})_{1 \leq k \leq K_0, 1 \leq h \leq H_0}, & \hat{\kappa}_{kh} &\equiv \frac{\frac{1}{|I_k| |J_h|} \sum_{i=1}^{|I_k|} \sum_{j=1}^{|J_h|} (A_{ij} - \hat{B}_{kh})^4}{\hat{S}_{kh}^4}. \tag{3.109}
\end{aligned}$$

As in the proof of Theorem 3.4.1, we consider the probability of the event \mathcal{F}_m that $[(n_k)_{1 \leq k \leq K}, (p_h)_{1 \leq h \leq H}, \tilde{\theta}, \tilde{\kappa}] = [(|I_k|)_{1 \leq k \leq K_0}, (|J_h|)_{1 \leq h \leq H_0}, \hat{\theta}, \hat{\kappa}]$ holds. The consistency assumption (vi) guarantees that if $(K_0, H_0) = (K, H)$, $\Pr(\mathcal{F}_m^c)$ converges to 0 in the limit

3. Statistical test on the number of biclusters in a latent block model

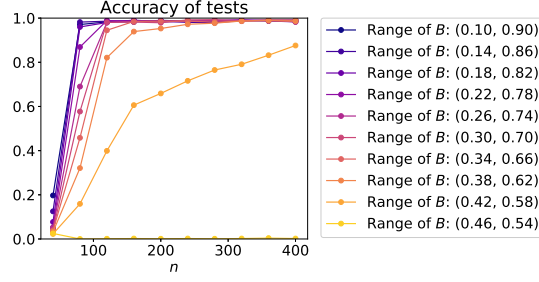


Figure 3.E1: Accuracy of the Jarque–Bera test under 10 different settings of block-wise mean B (Gaussian LBM).

of $m \rightarrow \infty$. Therefore, $\left|T^{\text{JB}} - \tilde{T}^{\text{JB}}\right|$ also converges in probability to zero, where

$$T^{\text{JB}} \equiv \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} T_{kh}^{\text{JB}},$$

$$T_{kh}^{\text{JB}} \equiv \frac{|I_k||J_h|}{6} \left[\hat{\theta}_{kh}^2 + \frac{1}{4}(\hat{\kappa}_{kh} - 3)^2 \right], \text{ for all } k = 1, \dots, K_0, h = 1, \dots, H_0. \quad (3.110)$$

From Slutsky’s theorem, $T^{\text{JB}} = \tilde{T}^{\text{JB}} + (T^{\text{JB}} - \tilde{T}^{\text{JB}})$ converges in law to the chi-squared distribution with $2K_0H_0$ degrees of freedom. Based on this discussion, to determine the number of blocks, we can use T^{JB} in (3.110) as a test statistic.

As in Section 3.5.3, we evaluated the accuracy of the above Jarque–Bera test. We generated data matrices from Gaussian LBMs with $(K, H) = (4, 3)$. We used the same experimental settings as in Section 3.5.3 except for the test method.

Figure 3.E1 shows the accuracy of the Jarque–Bera test under 10 different settings of block-wise mean B . By comparing Figures 6 and 3.E1, we see that the proposed test achieved comparable or better performance than the Jarque–Bera one in most settings (especially when the matrix size is small).

Currently, there is no other way than the proposed and Jarque–Bera tests for testing the cluster numbers of LBMs. A possible option might be to construct a test based on the Frobenius norm of matrix \hat{Z} instead of the operator one. However, as also pointed out in a study [92], such a test is expected to have lower power than the proposed test, since it does not aggressively detect the direction in which the matrix \hat{Z} deviates from the one generated from the null hypothesis, a distribution with zero mean and unit variance.

Chapter 4

Statistical test on the estimated bicluster structure of a relational data matrix

Model selection in latent block models has been a challenging but important task in the field of statistics. Specifically, a major challenge is encountered when constructing a test on a block structure obtained by applying a specific clustering algorithm to a finite size matrix. In this case, it becomes crucial to consider the selective bias in the block structure, that is, the block structure is selected from all the possible cluster memberships based on some criterion by the clustering algorithm. To cope with this problem, this chapter provides a selective inference method for latent block models. Specifically, we construct a statistical test on a set of row and column cluster memberships of a latent block model, which is given by a squared residue minimization algorithm. The proposed test, by its nature, includes and thus can also be used as the test on the set of row and column cluster numbers. We also propose an approximated version of the test based on simulated annealing to avoid combinatorial explosion in searching the optimal block structure. The results show that the proposed exact and approximated tests work effectively, compared to the naive test that does not take the selective bias into account.

4.1 Introduction

As we have described in Chapter 3, a latent block model or an LBM [58, 63] has been widely used as a generative model of a relational data matrix. In LBMs, we assume that there is an underlying block structure (i.e., a set of row and column cluster memberships) behind the observed data matrix and that each element of the matrix is generated independently from an identical distribution, given such a block structure. Particularly, a Gaussian LBM [104, 111] is useful to model a relational data matrix with real elements; this type of LBM is the focus of the current study. In a Gaussian LBM, we assume that each entry follows a Gaussian distribution, whose mean and variance are fixed constants in the same block (a

4. Statistical test on the estimated bicluster structure of a relational data matrix

formal description of Gaussian LBMs is given in Section 4.2.1).

Besides estimating the block structure from a given observed data matrix based on an LBM, it is also important to *test* the validity of a model (i.e., the number of blocks) or an estimation result. Until now, several tests have been proposed for determining the number of blocks in block models [16, 67, 92, 152, 164] (see Chapter 3 for details). Among these studies, only Chapter 3's test [152] can be applied to the LBM setting; however, its target is limited to the number of blocks, not to the cluster memberships. Moreover, it is an asymptotic test, and thus its guarantee cannot be verified with a finite size observed matrix.

In regard to an SBM, several studies have proposed a statistical test for a given set of community memberships of an observed matrix [53, 67, 74]. In [53], based on the numbers of edges within and across the clusters, two tests were proposed for an SBM; one of these tests included a goodness-of-fit test of community memberships. Although this study's objective is similar to ours, its problem setting is quite different from ours in various aspects, such as the setting of the alternative hypothesis and the assumptions in the network structure (e.g., there are two equal-sized communities in a given network and more intra-community edges than inter-community ones). Another study [67] proposed an asymptotic test on both the number of communities and the community memberships of an SBM, whose validity is guaranteed with the infinite matrix size. This study is different from ours in that our proposed test in this chapter is validated with a finite size matrix. [74] proposed a non-asymptotic test for an SBM setting; they generate finite samples of networks from the distribution of an SBM, conditioned on its sufficient statistics based on Markov chain Monte Carlo (MCMC), and, subsequently, compute the estimator of the p -value as the ratio of the test statistics of sampled networks being equal to or larger than that of an observed network. This study is somewhat similar to ours in that it tries to approximate the p -value under the condition that some function value of an observed matrix is given; however, it is fully based on a Metropolis-Hastings (MH) algorithm, and thus the resulting p -value is *not* exact with finite samples.

There have been many studies on statistical tests for SBMs, but none of them have enabled us to test the cluster memberships of LBMs. Particularly, in this chapter, we derive an *exact* p -value in the following context, which is a typical case in practice. First, we estimate an underlying block structure or cluster memberships of the rows and columns of an observed data matrix, based on a specific criterion. For instance, as a criterion, we use the *squared residue* or the sample variance within the same block [34, 63], whose formal definition is given in Section 4.2.2. Subsequently, we perform a statistical test on the clustering result, *which has been selected as an optimal block structure based on the data matrix*, in terms of the criterion described above. In regard to the construction of a valid statistical test, one concern is that it necessitates taking into account the *selective bias* [13, 88, 103]. A test on cluster memberships tends to be inappropriately positive, that is, it tends not to reject the hypothesis that the estimated cluster memberships are correct, when the test fails to consider the fact that the hypothetical set of cluster memberships was

4. Statistical test on the estimated bicluster structure of a relational data matrix

selected by using the information of a data matrix.

In order to perform a valid statistical inference in such a situation, [13, 88] introduced the methodology of post-selection inference. Particularly, *selective inference* methods facilitate inference of a hypothesis selected based on some criterion, where we use the same data for the hypothesis selection as well as for its inference [88]. The main idea behind the selective inference is to reveal the probability distribution of a given test statistic under the selection condition. By conditioning on the selection event, we can appropriately construct a test without the selective bias. Such selective inference methods have been developed for various problem settings, including variable selection in linear regression with L1 regularization [88] and that with marginal screening [90] and k-means clustering [69]. Concerning the problems related to the analysis of relational data matrices, several studies have proposed selective inference methods for biclustering [64, 89]. Although they also concern a block (or multiple blocks) in a relational data matrix, their problem settings are different from ours. In our problem setting, a block structure corresponds to a set of cluster memberships of *all* the rows and columns of an observed matrix. In other words, by rearranging the indices of rows and columns, a block structure is represented by a regular lattice on a matrix. However, [64, 89] aimed to find a submatrix (or multiple submatrices) of the original data matrix whose mean is significantly larger than zero. Figure 1 illustrates the difference between the optimal cluster memberships of the proposed and existing methods [89]¹. Since they are based on the mutually different assumptions on the latent bicluster structure, their “optimal” cluster memberships are not always identical, even with the same observed matrix. No study has proposed a selective inference method for the LBM setting, despite the effectiveness of LBMs in relational data analysis.

This chapter proposes a new selective inference method for LBMs. Unlike Chapter 3, where the validity of the test is guaranteed only in the asymptotic sense (i.e., with the infinite matrix size), we develop an **exact** test on a block structure, which is selected based on a given observed matrix with a **finite** size and the squared residue minimization algorithm. To construct such a statistical test, we consider the fact that the selection event based on the squared residue can be formulated as a set of quadratic inequalities in terms of the data vector, which is the vectorization of the observed data matrix. On this basis, we can show that the test statistic follows a truncated chi distribution, under the selection condition (a formal definition of the test statistic is given in Section 4.3).

¹To plot Figure 1, we independently generated data matrices with the sizes of $(n, p) = (9, 9)$. We set the null and hypothetical sets of cluster numbers at $(2, 2)$; we defined the null cluster memberships as $g_i^{(N),(1)} = (i \bmod 2) + 1$, for all i , and $g_j^{(N),(2)} = (j \bmod 2) + 1$ for all j . In regard to the mean vector, we used the following setting: $\boldsymbol{\mu}_0 = \text{vec} \left(\begin{bmatrix} 0.5 & 0 \\ 0 & 0 \end{bmatrix} \right)$. Based on the above settings, we generated a data vector by $\boldsymbol{x} \sim N(\boldsymbol{\mu}_0, 0.75^2 I)$ and applied the biclustering algorithms of the proposed and existing methods [89]. The biclustering algorithm of the proposed method outputs a regular-grid bicluster structure based on the squared residue minimization, while that of [89] outputs an $n_0 \times p_0$ submatrix with the largest sample mean, where we set $n_0 = p_0 = 5$.

4. Statistical test on the estimated bicluster structure of a relational data matrix

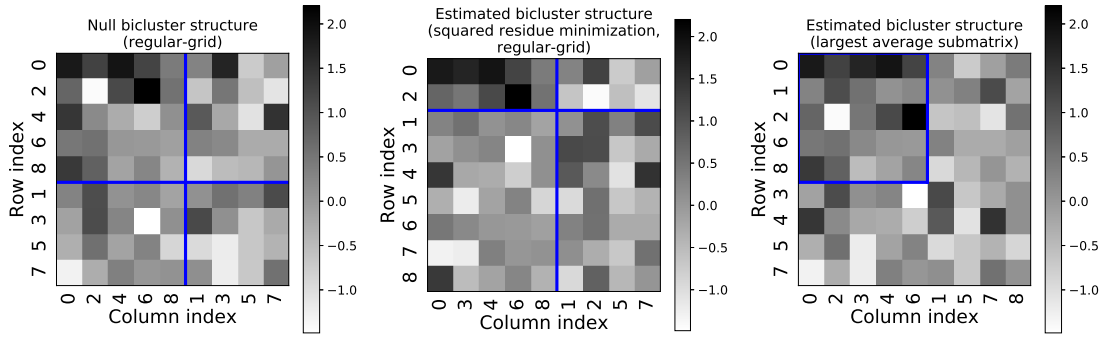


Figure 1: Examples of the null and estimated bicluster structures of the observed data matrix with the size of $(n, p) = (9, 9)$. The rows and columns of the observed matrix were sorted according to their clusters, and the blue lines indicate the cluster memberships. The biclustering algorithms of the proposed and existing methods [89] do not necessarily yield identical bicluster structures with the same observed matrix.

Since the exact test requires solving two combinatorial optimization problems—one for selecting the block structure with the minimum squared residue, and the other for determining the truncation interval of the distribution of the test statistic—its computation will be intractable with a large size observed matrix or with a larger hypothetical number of blocks. To cope with such combinatorial explosion, we also develop an approximated version of the test based on simulated annealing (SA).

The remaining part of this chapter is as follows. In Section 4.2, we first define notations and describe assumptions necessary for developing the proposed statistical test. We also define the squared residue, which we use for measuring the quality of a given set of row and column cluster memberships. In Section 4.3, we give the formal statement of the null and alternative hypotheses of the proposed test, define the test statistic, and derive its null distribution. Our main contribution lies in Theorem 4.3.1; it states that, under the null hypothesis, the test statistic follows a truncated chi distribution, whose truncation interval is determined by the selection result. We also give an approximated version of the test. In Section 4.4, we experimentally show the effectiveness of the proposed exact and approximated tests, by checking the behavior of the p -values and measuring the true and false positive ratios (TPR and FPR) in both the realizable [i.e., the hypothetical cluster numbers of rows and columns (K_0, H_0) are equal to the null ones (K, H)] and unrealizable (i.e., at least one of $K_0 < K$ and $H_0 < H$ holds) cases. Finally, we discuss the findings and conclude the chapter in Section 4.5 and Section 4.6, respectively.

4.2 Problem settings

4.2.1 Notations and assumptions on data matrix

Throughout this chapter, we use the following definitions and notations.

- Let $A = (A_{ij})_{1 \leq i \leq n, 1 \leq j \leq p} \in \mathbb{R}^{n \times p}$ be an observed data matrix with the size of $n \times p$. When constructing a statistical test, it is more convenient to use the vector representation of matrix A , instead of A itself:

$$\mathbf{x} = \text{vec}(A) \in \mathbb{R}^{np}, \quad x_{n(j-1)+i} = A_{ij}, \quad \text{for } i = 1, \dots, n, \quad j = 1, \dots, p. \quad (4.1)$$

- Let $g_i^{(1)}$ be the cluster index of the i th row, and $g^{(1)} = (g_i^{(1)})_{1 \leq i \leq n}$. Similarly, let $g_j^{(2)}$ be the cluster index of the j th column, and $g^{(2)} = (g_j^{(2)})_{1 \leq j \leq p}$. We denote a set of row and column clusters as $g = (g^{(1)}, g^{(2)}) \in \mathcal{G}$, where $\mathcal{G} = \{(g^{(1)}, g^{(2)})\}$ is a set of all possible cluster memberships. We also define that $\mathcal{G}_{K_0 H_0}$ is a set of all possible cluster memberships with $K_0 \times H_0$ or less blocks.
- In the null hypothesis of the proposed test, we assume that there exists a set of block memberships $g^{(N)} = (g^{(N),(1)}, g^{(N),(2)})$ and that, given $g^{(N)}$, each (i, j) th element A_{ij} of an observed matrix A is generated independently from a Gaussian distribution with a block-wise (**unknown**) mean $P_{ij} \equiv B_{g_i^{(N),(1)} g_j^{(N),(2)}}$ and (**known**²) variance σ_0^2 , where B_{kh} is the mean of the (k, h) th block:

$$A_{ij} \sim N(P_{ij}, \sigma_0^2), \quad \text{for all } i = 1, \dots, n, \quad j = 1, \dots, p. \quad (4.2)$$

In vector representation, this assumption is given by

$$\mathbf{x} \sim N(\boldsymbol{\mu}_0, \sigma_0^2 I), \quad (4.3)$$

where $\boldsymbol{\mu}_0$ is the **unknown** block-wise mean vector.

- Let (K, H) be the minimum set of row and column cluster numbers required to represent the above null set of block memberships $g^{(N)}$. In the proposed test, we fix a hypothetical set of cluster numbers (K_0, H_0) , estimate the block structure of an observed matrix with $K_0 \times H_0$ blocks, and perform a test on the estimated block memberships, which, by its nature, includes a test on cluster numbers (i.e., $(K, H) = (K_0, H_0)$ or at least one of $K_0 < K$ and $H_0 < H$ holds)³.

²We also derive the null distribution of a test statistic in case that variance σ_0^2 is unknown in Appendix 4.G.

³It must be noted, however, that the proposed test cannot be applied directly for sequential testing on cluster numbers, where the hypothetical numbers of clusters are tested in ascending order (i.e., $(K_0, H_0) = (1, 1), (1, 2), (2, 1), \dots$) until the null hypothesis is accepted. This is because the proposed test cannot distinguish the following two alternative cases: (1) $(K, H) = (K_0, H_0)$ holds, however, the estimated cluster memberships are incorrect, and (2) $K_0 < K$ or $H_0 < H$ holds.

4. Statistical test on the estimated bicluster structure of a relational data matrix

- We denote the set of rows in the k th cluster as $I_k = \{i : g_i^{(1)} = k\}$. Similarly, we denote the set of columns in the h th cluster as $J_h = \{j : g_j^{(2)} = h\}$.
- We denote the cluster membership vector of rows as follows:

$$\bar{\mathbf{e}}^{(k)} = (\bar{e}_i^{(k)})_{1 \leq i \leq n} \in \mathbb{R}^n, \bar{e}_i^{(k)} = \begin{cases} \frac{1}{\sqrt{|I_k|}} & \text{if } g_i^{(1)} = k, \\ 0 & \text{otherwise.} \end{cases} \quad (4.4)$$

Similarly, we denote the cluster membership vector of columns as follows:

$$\underline{\mathbf{e}}^{(h)} = (\underline{e}_j^{(h)})_{1 \leq j \leq p} \in \mathbb{R}^p, \underline{e}_j^{(h)} = \begin{cases} \frac{1}{\sqrt{|J_h|}} & \text{if } g_j^{(2)} = h, \\ 0 & \text{otherwise.} \end{cases} \quad (4.5)$$

Based on these vectors $\bar{\mathbf{e}}^{(k)}$ and $\underline{\mathbf{e}}^{(h)}$, we define a vector $\mathbf{e}^{(k,h)} \equiv \underline{\mathbf{e}}^{(h)} \otimes \bar{\mathbf{e}}^{(k)} \in \mathbb{R}^{np}$ and matrix $E^{(g)} \equiv I - \sum_k \sum_h \mathbf{e}^{(k,h)} (\mathbf{e}^{(k,h)})^\top$. It must be noted that $E^{(g)}$ is a projection matrix, that is, $(E^{(g)})^\top = E^{(g)}$ and $(E^{(g)})^2 = E^{(g)}$ hold.

4.2.2 Clustering algorithm based on squared residue minimization

To estimate the block structure of a given observed matrix A , we use a clustering algorithm $\mathcal{A} : \mathbf{x} \mapsto \hat{\mathcal{M}} \in \mathcal{G}_{K_0 H_0}$ that outputs a block structure minimizing the *squared residue*, that is, the sample variance σ^2 within the same block. A squared residue has been proposed for measuring the quality of a biclustering result [34, 63], and its definition is given by

$$\begin{aligned} \sigma^2 &= \frac{1}{np} \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \sum_{i \in I_k} \sum_{j \in J_h} \left(A_{ij} - \frac{1}{|I_k||J_h|} \sum_{i' \in I_k} \sum_{j' \in J_h} A_{i'j'} \right)^2 \\ &= \frac{1}{np} \left[\sum_{i,j} A_{ij}^2 - \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \frac{1}{|I_k||J_h|} \left(\sum_{i \in I_k} \sum_{j \in J_h} A_{ij} \right)^2 \right] \\ &= \frac{1}{np} \left[\sum_{i,j} A_{ij}^2 - \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \left(\frac{1}{\sqrt{|I_k||J_h|}} \sum_{i \in I_k} \sum_{j \in J_h} A_{ij} \right)^2 \right] \\ &= \frac{1}{np} \left\{ \mathbf{x}^\top \mathbf{x} - \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} [(\underline{\mathbf{e}}^{(h)} \otimes \bar{\mathbf{e}}^{(k)})^\top \mathbf{x}]^2 \right\} \\ &= \frac{1}{np} \left\{ \mathbf{x}^\top \mathbf{x} - \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} [(\mathbf{e}^{(k,h)})^\top \mathbf{x}]^2 \right\} \\ &= \frac{1}{np} \mathbf{x}^\top \left[I - \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \mathbf{e}^{(k,h)} (\mathbf{e}^{(k,h)})^\top \right] \mathbf{x} = \frac{1}{np} \mathbf{x}^\top E^{(g)} \mathbf{x}. \end{aligned} \quad (4.6)$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

Therefore, the squared residue minimization clustering algorithm \mathcal{A} outputs the set of cluster memberships $\hat{g} = (\hat{g}^{(1)}, \hat{g}^{(2)})$, which satisfies

$$\hat{g} \in \hat{\mathcal{M}}(\mathbf{x}) = \arg \min_{g \in \mathcal{G}_{K_0 H_0}} \sigma^2 = \arg \min_{g \in \mathcal{G}_{K_0 H_0}} \mathbf{x}^\top E^{(g)} \mathbf{x}. \quad (4.7)$$

It must be noted that the above solution \hat{g} is the maximum likelihood estimator of the cluster memberships with a mean estimator $\hat{B}(g)$ and a known standard deviation σ_0 . The log likelihood of a set of cluster memberships $g = (g^{(1)}, g^{(2)})$ and the mean parameter B is given by

$$\begin{aligned} \mathcal{L}(g, B; \mathbf{x}) &= -np \log \left(\sqrt{2\pi\sigma_0^2} \right) - \frac{1}{2\sigma_0^2} \sum_{i=1}^n \sum_{j=1}^p \left(x_{n(j-1)+i} - B_{g_i^{(1)} g_j^{(2)}} \right)^2 \\ &= -np \log \left(\sqrt{2\pi\sigma_0^2} \right) - \frac{1}{2\sigma_0^2} \sum_{k=1}^{K_0} \sum_{h=1}^{H_0} \sum_{i \in I_k} \sum_{j \in J_h} \left(x_{n(j-1)+i} - B_{kh} \right)^2. \end{aligned} \quad (4.8)$$

Let $\hat{B}(g) = (\hat{B}_{kh}(g))_{1 \leq k \leq K_0, 1 \leq h \leq H_0}$ be the maximum likelihood estimator of mean B for a given fixed cluster memberships g . From (4.8), we can easily derive that $\hat{B}_{kh}(g) = (1/|I_k||J_h|) \sum_{i \in I_k} \sum_{j \in J_h} x_{n(j-1)+i}$. By combining this fact with (4.6) and (4.8), we see that the squared residue minimization is equivalent to the likelihood maximization with a mean estimator $\hat{B}(g)$.

Equation (4.7) is equivalent to a set of quadratic inequalities

$$\mathbf{x}^\top E^{(\hat{g})} \mathbf{x} \leq \mathbf{x}^\top E^{(g)} \mathbf{x} \iff \mathbf{x}^\top (E^{(g)} - E^{(\hat{g})}) \mathbf{x} \geq 0, \quad (4.9)$$

for all $g \in \mathcal{G}_{K_0 H_0}$. In other words, the selection rule can be represented as a set of quadratic inequalities in terms of the data vector \mathbf{x} . It must be noted that, under the null hypothesis, the solution \hat{g} of (4.7) is unique almost surely. To prove this fact, we first define a quadratic function $F^{(g, g')} : \mathbb{R}^{np} \mapsto \mathbb{R}$ for a fixed (g, g') as $F^{(g, g')}(\mathbf{x}) \equiv \mathbf{x}^\top (E^{(g)} - E^{(g')}) \mathbf{x}$. We also define that $g = g'$, if the sets of cluster memberships g and g' are equivalent up to the permutation of cluster indices, and that $g \neq g'$ otherwise. If $g \neq g'$, we have

$$E^{(g)} - E^{(g')} \neq 0 \quad (\because \text{the proof is in Appendix 4.A}), \quad (4.10)$$

and thus the Lebesgue measure of a set of points \mathbf{x} that satisfy $F^{(g, g')}(\mathbf{x}) = 0$ is zero. By combining this fact and the assumption (4.3) of the null hypothesis, $F^{(g, g')}(\mathbf{x}) \neq 0$ holds for a fixed combination of (g, g') almost surely. Since $\{(g, g') : g, g' \in \mathcal{G}_{K_0 H_0}, g \neq g'\}$ is a finite set, we finally have

$$\Pr \left(\exists g, g' \in \mathcal{G}_{K_0 H_0}, \text{ s.t. } g \neq g', g, g' \in \arg \min_{g \in \mathcal{G}_{K_0 H_0}} \sigma^2 \right)$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

$$\begin{aligned}
&= \Pr \left(\exists g, g' \in \mathcal{G}_{K_0 H_0}, \text{ s.t. } g \neq g', F^{(g, g')}(\mathbf{x}) = 0, g, g' \in \arg \min_{g \in \mathcal{G}_{K_0 H_0}} \sigma^2 \right) \\
&\leq \Pr \left(\exists g, g' \in \mathcal{G}_{K_0 H_0}, \text{ s.t. } g \neq g', F^{(g, g')}(\mathbf{x}) = 0 \right) \\
&\leq \sum_{g, g' \in \mathcal{G}_{K_0 H_0}, g \neq g'} \Pr \left(F^{(g, g')}(\mathbf{x}) = 0 \right) = 0.
\end{aligned} \tag{4.11}$$

In case of a tie (i.e., multiple solutions of \hat{g} exist that satisfy (4.7)) that occurs with probability zero, we can choose any one of them as \hat{g} independently with \mathbf{x} .

4.3 Main results: Statistical test on the solution of squared residue minimization

4.3.1 Null distribution of test statistic T

As described in Section 4.2, in the null hypothesis of the proposed test, we assume that there exists a set of block memberships $g^{(N)}$ and that given $g^{(N)}$, each element of an observed data vector \mathbf{x} is generated independently from a Gaussian distribution, whose mean is constant within the same block. Our main purpose is to test whether an estimated block structure $\hat{g} \equiv (\hat{g}^{(1)}, \hat{g}^{(2)})$, which is selected based on the squared residue criterion in Section 4.2.2, is equal to the null one $g^{(N)}$. Formally, the null and alternative hypotheses of the proposed test are given by

$$(N) : E^{(\hat{g})} \boldsymbol{\mu}_0 = \mathbf{0}, \quad (A) : E^{(\hat{g})} \boldsymbol{\mu}_0 \neq \mathbf{0}. \tag{4.12}$$

It must be noted that the equation $E^{(\hat{g})} \boldsymbol{\mu}_0 = \mathbf{0}$ is equivalent to the statement that the elements of the vector $\boldsymbol{\mu}_0$ are constant in the same block in the set of cluster memberships \hat{g} . In other words, the above statement of the null hypothesis is that a given observed matrix is generated based on the latent block structure \hat{g} , which is selected as a solution that minimizes the squared residue.

To perform the test of (4.12), we check the squared residue σ^2 of the given observed matrix A under the condition that the estimated block structure \hat{g} is selected. Under the null hypothesis, we have $E^{(\hat{g})} \boldsymbol{\mu}_0 = \mathbf{0}$. Here, matrix $E^{(\hat{g})}$ solely depends on the estimated set of cluster memberships \hat{g} . In other words, under the condition that $\hat{\mathcal{M}}(\mathbf{x}) = \hat{g}$ holds, matrix $E^{(\hat{g})}$ is fixed. Therefore, based on the result in [103], the following theorem holds:

Theorem 4.3.1. *Under the null hypothesis, we have*

$$T \equiv \frac{\|\mathbf{r}\|_2}{\sigma_0}, \quad T | \{\hat{g}, \mathbf{z}, \mathbf{u}\} \sim \chi_{(np - K_0 H_0) | \hat{M}(\hat{g})}, \tag{4.13}$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

where $\|\cdot\|_2$ and $\chi_{c|M}$, respectively, denote the Euclid norm and the truncated chi distribution with c degrees of freedom and with truncation interval of M and

$$\begin{aligned} \mathbf{r} &\equiv E^{(\hat{g})}\mathbf{x}, \quad \mathbf{u} \equiv \frac{1}{\|\mathbf{r}\|_2}\mathbf{r}, \quad \mathbf{z} \equiv \mathbf{x} - \mathbf{r}, \\ \hat{M}^{(\hat{g})} &\equiv \{t \geq 0 : \hat{g} \in \hat{\mathcal{M}}(t\sigma_0\mathbf{u} + \mathbf{z})\}. \end{aligned} \quad (4.14)$$

Proof. Let E be a fixed $np \times np$ projection matrix satisfying the following conditions:

- $\text{rank}(E) = np - K_0H_0$.
- $E\boldsymbol{\mu}_0 = \mathbf{0}$.

A singular value decomposition of a matrix E satisfying the above two conditions is given by

$$E = V^\top DV, \quad D \equiv \begin{bmatrix} I_{(np-K_0H_0)} & O_{(np-K_0H_0),K_0H_0} \\ O_{K_0H_0,(np-K_0H_0)} & O_{K_0H_0,K_0H_0} \end{bmatrix}, \quad (4.15)$$

where we denote the $a \times a$ identity matrix and $a \times b$ zero matrix, respectively, as I_a and $O_{a,b}$.

Based on such a matrix E , we use the following notations:

$$\mathbf{r}_E \equiv E\mathbf{x}, \quad T_E = \frac{\|\mathbf{r}_E\|_2}{\sigma_0}, \quad \mathbf{u}_E \equiv \frac{1}{\|\mathbf{r}_E\|_2}\mathbf{r}_E, \quad \mathbf{z}_E \equiv \mathbf{x} - \mathbf{r}_E. \quad (4.16)$$

We can transform T_E by the following equations:

$$\begin{aligned} T_E &= \frac{\sqrt{\mathbf{x}^\top E\mathbf{x}}}{\sigma_0} = \frac{\sqrt{(\mathbf{x} - \boldsymbol{\mu}_0)^\top E(\mathbf{x} - \boldsymbol{\mu}_0)}}{\sigma_0} \quad (\because E\boldsymbol{\mu}_0 = \mathbf{0}) \\ &= \frac{\sqrt{(\mathbf{x} - \boldsymbol{\mu}_0)^\top V^\top DV(\mathbf{x} - \boldsymbol{\mu}_0)}}{\sigma_0} = \frac{\sqrt{(\mathbf{x} - \boldsymbol{\mu}_0)^\top V^\top \tilde{D}^\top \tilde{D}V(\mathbf{x} - \boldsymbol{\mu}_0)}}{\sigma_0}, \\ \tilde{D} &\equiv \begin{bmatrix} I_{(np-K_0H_0)} & O_{(np-K_0H_0),K_0H_0} \end{bmatrix} \in \mathbb{R}^{(np-K_0H_0) \times np}. \end{aligned} \quad (4.17)$$

Here, we used the fact that $\tilde{D}^\top \tilde{D} = D$.

By using the assumption that $\mathbf{x} \sim N(\boldsymbol{\mu}_0, \sigma_0^2 I)$ holds and the independence of matrix E of \mathbf{x} , we have

$$\begin{aligned} \frac{1}{\sigma_0} \tilde{D}V(\mathbf{x} - \boldsymbol{\mu}_0) &\sim N(\mathbf{0}, \tilde{D}V(\tilde{D}V)^\top) \\ \iff \frac{1}{\sigma_0} \tilde{D}V(\mathbf{x} - \boldsymbol{\mu}_0) &\sim N(\mathbf{0}, I_{np-K_0H_0}). \end{aligned} \quad (4.18)$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

Here, we used the fact that $\tilde{D}\tilde{D}^\top = I$. Therefore, by combining (4.17) and (4.18), we have

$$T_E \sim \chi_{(np-K_0H_0)}, \quad (4.19)$$

where χ_c denotes the chi distribution with c degrees of freedom.

In regard to \mathbf{u}_E and \mathbf{z}_E , we have

$$\mathbf{u}_E \cdot \mathbf{z}_E = \frac{1}{\|\mathbf{r}_E\|_2} \mathbf{r}_E^\top (\mathbf{x} - \mathbf{r}_E) = \frac{1}{\|\mathbf{r}_E\|_2} (\mathbf{x}^\top E^\top \mathbf{x} - \mathbf{x}^\top E^\top E \mathbf{x}) = 0. \quad (4.20)$$

In the last equation, we considered the fact that $E^\top E = E$.

Here, since T_E and $(\mathbf{u}_E, \mathbf{z}_E)$ are mutually independent (the proof of this is in Appendix 4.D), we have

$$T_E | \mathbf{u}_E, \mathbf{z}_E \sim \chi_{(np-K_0H_0)}. \quad (4.21)$$

Next, we consider adding a condition of selection event of \hat{g} to the distribution of $T_E | \mathbf{u}_E, \mathbf{z}_E$ in (4.21). Given \mathbf{u}_E and \mathbf{z}_E , the result of selection depends solely on the value of T_E . Therefore, adding the selection condition $\hat{\mathcal{M}}(\mathbf{u}_E T_E \sigma_0 + \mathbf{z}_E) = \hat{g}$ to (4.21) corresponds to truncation of T_E to the region where $\hat{\mathcal{M}}(\mathbf{u}_E T_E \sigma_0 + \mathbf{z}_E) = \hat{g}$ holds:

$$T_E | \mathbf{u}_E, \mathbf{z}_E, \hat{g} \sim \chi_{(np-K_0H_0) | \hat{\mathcal{M}}(\hat{g})(E)}. \quad (4.22)$$

Third, we consider replacing E in (4.22) with $E^{(\hat{g})}$, which is the output by clustering algorithm \mathcal{A} based on the data vector \mathbf{x} . It must be noted that the matrix $E^{(\hat{g})}$ is also a projection matrix that satisfies

$$\text{rank}(E^{(\hat{g})}) = np - K_0H_0 \quad (\because \text{the proof is in Appendix 4.B}), \quad (4.23)$$

from its definition, and $E^{(\hat{g})} \boldsymbol{\mu}_0 = \mathbf{0}$ holds.

Since the matrix $E^{(\hat{g})}$ depends on the data vector \mathbf{x} only through the choice of \hat{g} (i.e., $E^{(\hat{g})}$ is fixed, given \hat{g}), under the condition that the selection result \hat{g} is given, (4.22) still holds with $E^{(\hat{g})}$, which concludes the proof. \square

Remark 4.3.1 (Generalization of Theorem 4.3.1). *Theorem 4.3.1 holds if the selection event of the estimated block structure \hat{g} can be formulated as a set of quadratic inequalities in terms of the data vector \mathbf{x} , by modifying the definition of the function \mathcal{M} . In other words, for a selected block structure \hat{g} , there exists some $\mathcal{I}_{\hat{g}} \in \mathbb{N}$ and $\{Q^{(\hat{g},i)}, \boldsymbol{\alpha}^{(\hat{g},i)}, \beta^{(\hat{g},i)}\}$, $i = 1, \dots, \mathcal{I}_{\hat{g}}$, and the selection event of \hat{g} is represented by*

$$\hat{g} \in \mathcal{M}(\mathbf{x}) \iff \bigcap_{i \in \{1, \dots, \mathcal{I}_{\hat{g}}\}} \{ \mathbf{x}^\top Q^{(\hat{g},i)} \mathbf{x} + (\boldsymbol{\alpha}^{(\hat{g},i)})^\top \mathbf{x} + \beta^{(\hat{g},i)} \geq 0 \}. \quad (4.24)$$

Let $g(i)$ be the i th pattern of all the block structures with $K_0 \times H_0$ blocks or less, where $i = 1, \dots, |\mathcal{G}_{K_0H_0}|$. Then, if we set $\mathcal{I}_{\hat{g}} \equiv |\mathcal{G}_{K_0H_0}|$, $Q^{(\hat{g},i)} = E^{(g(i))} - E^{(\hat{g})}$, $\boldsymbol{\alpha}^{(\hat{g},i)} = \mathbf{0}$, and $\beta^{(\hat{g},i)} = 0$, the selection event in (4.24) will lead to the use of a squared residue solution.

4. Statistical test on the estimated bicluster structure of a relational data matrix

It must be noted that if there exists multiple sets of cluster memberships that minimize the squared residue σ^2 , which occurs with probability zero, from the discussion in Section 4.2.2, then Theorem 4.3.1 will hold for any one of them. Moreover, we define that a set of block memberships g' is a *refinement* of g iff any block in g' is a submatrix of some block in g . If \hat{g}' is a refinement of \hat{g} , then Theorem 4.3.1 will also hold when \hat{g} is replaced by \hat{g}' . In other words, we cannot detect that a given block structure represents a “finer division than necessary” with the proposed test; solving this problem is beyond the scope of this dissertation.

4.3.2 Statistical test based on truncated chi distribution

To perform a statistical test based on Theorem 4.3.1, we have to specify the truncation interval of $\hat{M}^{(\hat{g})} \equiv \{t \geq 0 : \hat{\mathcal{M}}(t\sigma_0\mathbf{u} + \mathbf{z}) = \hat{g}\}$. As shown in (4.9), this is equivalent to an interval satisfying the following condition for all g :

$$\begin{aligned} (t\sigma_0\mathbf{u} + \mathbf{z})^\top (E^{(g)} - E^{(\hat{g})}) (t\sigma_0\mathbf{u} + \mathbf{z}) &\geq 0 \\ \iff f^{(g,\hat{g})}(t) \equiv a^{(g,\hat{g})}t^2 + b^{(g,\hat{g})}t + c^{(g,\hat{g})} &\geq 0, \end{aligned} \quad (4.25)$$

where

$$\begin{aligned} a^{(g,\hat{g})} &\equiv \sigma_0^2 \mathbf{u}^\top (E^{(g)} - E^{(\hat{g})}) \mathbf{u}, \\ b^{(g,\hat{g})} &\equiv \sigma_0 [\mathbf{u}^\top (E^{(g)} - E^{(\hat{g})}) \mathbf{z} + \mathbf{z}^\top (E^{(g)} - E^{(\hat{g})}) \mathbf{u}], \\ c^{(g,\hat{g})} &\equiv \mathbf{z}^\top (E^{(g)} - E^{(\hat{g})}) \mathbf{z}. \end{aligned} \quad (4.26)$$

From the definition of \mathbf{u} and \mathbf{z} in (4.14), we have $E^{(\hat{g})}\mathbf{u} = \mathbf{u}$ and $E^{(\hat{g})}\mathbf{z} = \mathbf{0}$, which simplify the above coefficients $a^{(g,\hat{g})}$, $b^{(g,\hat{g})}$, and $c^{(g,\hat{g})}$ as follows:

$$\begin{aligned} a^{(g,\hat{g})} &= -\sigma_0^2 \mathbf{u}^\top (I - E^{(g)}) \mathbf{u} = -\sigma_0^2 \|(I - E^{(g)}) \mathbf{u}\|_2^2 \leq 0, \\ b^{(g,\hat{g})} &= 2\sigma_0 \mathbf{u}^\top E^{(g)} \mathbf{z}, \\ c^{(g,\hat{g})} &= \mathbf{z}^\top E^{(g)} \mathbf{z} = \|E^{(g)} \mathbf{z}\|_2^2 \geq 0. \end{aligned} \quad (4.27)$$

Here, in the transformation of $b^{(g,\hat{g})}$, we used the fact that matrices $E^{(g)}$ and $E^{(\hat{g})}$ are symmetric.

We consider the condition under which (4.25) holds in the two cases, $a^{(g,\hat{g})} = 0$ and $a^{(g,\hat{g})} < 0$.

- If $a^{(g,\hat{g})} = 0$, we have $E^{(g)}\mathbf{u} = \mathbf{u}$, which results in that $b^{(g,\hat{g})} = 2\sigma_0 \mathbf{u}^\top \mathbf{z} = 0$ (since $\mathbf{u}^\top \mathbf{z} = 0$ holds from (4.20)). Therefore, in this case, the selection condition (4.25) always holds.

4. Statistical test on the estimated bicluster structure of a relational data matrix

- If $a^{(g,\hat{g})} < 0$, $\max_t f^{(g,\hat{g})}(t) \geq f^{(g,\hat{g})}(0) = c^{(g,\hat{g})} \geq 0$. Therefore, for $t \geq 0$, the interval that satisfies $f^{(g,\hat{g})}(t) \geq 0$ is $\left[0, \frac{-b^{(g,\hat{g})} - \sqrt{(b^{(g,\hat{g})})^2 - 4a^{(g,\hat{g})}c^{(g,\hat{g})}}{2a^{(g,\hat{g})}}\right]$.

Overall, the interval of t where (4.25) holds is given by

$$\hat{M}^{(\hat{g})} = [0, \beta^{(\hat{g})}], \quad \beta^{(\hat{g})} \equiv \min_{g:a^{(g,\hat{g})} \neq 0} \left(\frac{-b^{(g,\hat{g})} - \sqrt{(b^{(g,\hat{g})})^2 - 4a^{(g,\hat{g})}c^{(g,\hat{g})}}{2a^{(g,\hat{g})}} \right). \quad (4.28)$$

It must be noted that $\bigcap_{g \in \mathcal{G}_{K_0 H_0}, g \neq \hat{g}} (a^{(g,\hat{g})} < 0)$ holds almost surely, based on a similar discussion as that in Section 4.2.2. Formally, for a fixed $g, g' \in \mathcal{G}_{K_0 H_0}$, $\mathbf{y} \equiv E^{(g')} \mathbf{x}$ follows a Gaussian distribution. If $g \neq g'$, $E^{(g)} - E^{(g')}$ is not a zero matrix, and thus the Lebesgue measure of a set of points \mathbf{y} satisfying $\mathbf{y}^\top (E^{(g)} - E^{(g')}) \mathbf{y} = 0$ is zero. Similarly, $\|\mathbf{y}\|_2^2 > 0$ holds with probability one. By combining these facts, $a^{(g,g')} \equiv \frac{1}{\|\mathbf{y}\|_2^2} \mathbf{y}^\top (E^{(g)} - E^{(g')}) \mathbf{y} \neq 0$ holds for a fixed combination of (g, g') satisfying $g \neq g'$ almost surely. Therefore, we have

$$\begin{aligned} & \Pr(\exists g \in \mathcal{G}_{K_0 H_0}, \text{ s.t. } g \neq \hat{g}, a^{(g,\hat{g})} = 0) \\ & \leq \Pr(\exists g, g' \in \mathcal{G}_{K_0 H_0}, \text{ s.t. } g \neq g', a^{(g,g')} = 0) \\ & \leq \sum_{g, g' \in \mathcal{G}_{K_0 H_0}, g \neq g'} \Pr(a^{(g,g')} = 0) = 0. \end{aligned} \quad (4.29)$$

To derive the last equation, we used the fact that $\{(g, g') : g, g' \in \mathcal{G}_{K_0 H_0}, g \neq g'\}$ is a finite set.

We denote a set of cluster memberships attaining the boundary of this interval as \tilde{g} , that is,

$$\tilde{g} \equiv \arg \min_{g:a^{(g,\hat{g})} \neq 0} \left(\frac{-b^{(g,\hat{g})} - \sqrt{(b^{(g,\hat{g})})^2 - 4a^{(g,\hat{g})}c^{(g,\hat{g})}}{2a^{(g,\hat{g})}} \right). \quad (4.30)$$

From Theorem 4.3.1, given $\{\hat{g}, \mathbf{z}, \mathbf{u}\}$, a p -value p_T of the test statistic T in (4.13) is given by

$$p_T = \begin{cases} 1 - \frac{\gamma\left(\frac{np - K_0 H_0}{2}, \frac{T^2}{2}\right)}{\gamma\left(\frac{np - K_0 H_0}{2}, \frac{(\beta^{(\hat{g})})^2}{2}\right)} \sim U[0, 1] & \text{if } 0 \leq T \leq \beta^{(\hat{g})}, \\ 0 & \text{otherwise,} \end{cases} \quad (4.31)$$

where $\gamma(\cdot, \cdot)$ is the lower incomplete gamma function. This holds from the fact that, for any random variable X with a probability density function $f(x)$, $F(X) \equiv \int_{-\infty}^X f(x) dx \sim$

4. Statistical test on the estimated bicluster structure of a relational data matrix

$U[0, 1]$. To derive the p -value in (4.31), we used the fact that the cumulative distribution function of the chi-square distribution with c degrees of freedom and with truncation interval of $[0, a]$ is given by

$$\begin{cases} F(x) = 0 & \text{if } x < 0, \\ F(x) = \frac{\gamma(c/2, x/2)}{\gamma(c/2, a/2)} & \text{if } 0 \leq x \leq a, \\ F(x) = 1 & \text{if } x > 1. \end{cases} \quad (4.32)$$

4.3.3 Approximated test based on simulated annealing

The exact statistical test in Section 4.3.2 requires us to find (i) the optimal set of cluster memberships \hat{g} , which minimizes the squared residue in (4.6), and (ii) the set of cluster memberships \tilde{g} in (4.30), which determines the truncation interval. We can see that the number of mutually different patterns of block structures with **exactly** $K_0 \times H_0$ blocks is lower bounded by $K_0^{n-K_0} H_0^{p-H_0}$ (see Appendix 4.C for more detailed discussions).

To cope with such combinatorial explosion, we propose an approximated statistical test based on SA, besides the exact test described in Section 4.3.2. SA is an iterative algorithm that can be used for obtaining approximated solutions of combinatorial optimization problems [79, 147]; its basic procedure is given as follows:

1. Define a cooling schedule or the sequence of temperatures $\{T_t\}_{t=0}^{\infty}$, a threshold ϵ , a finite set of states \mathcal{S} , and an objective function f on \mathcal{S} . For all the experiments, we set the threshold at $\epsilon = 10^{-6}$. Our purpose is to find a state $x \in \mathcal{S}$ that minimizes $f(x)$. For each state $x \in \mathcal{S}$, we also define a set of neighbors $N(x) \subseteq \mathcal{S}$ and a transition probability $R(x, x')$ from state x to x' , for all $x' \in \mathcal{S}$, where $R(x, x') > 0$ if $x' \in N(x)$ and $R(x, x') = 0$ otherwise. Finally, define an initial step $t \leftarrow 0$ and initial state $x_0 \in \mathcal{S}$, and let $f^{(0)} \equiv f(x_0)$.
2. If $T_t < \epsilon$, stop the algorithm and output the current state x_t . Otherwise, randomly choose a neighbor x' of the current state x_t (i.e., $x' \in N(x_t)$) with probability $R(x_t, x')$. Let $f' \equiv f(x')$ and $\Delta f \equiv f' - f^{(t)}$.
 - If $\Delta f < 0$, then move to state x' and set $x_{t+1} = x'$ and $f^{(t+1)} = f'$.
 - Otherwise, with probability $\exp\left(-\frac{\Delta f}{T_t}\right)$, move to state x' and set $x_{t+1} = x'$ and $f^{(t+1)} = f'$. Otherwise, stay at the current state x_t and set $x_{t+1} = x_t$ and $f^{(t+1)} = f^{(t)}$.
3. Let $t \leftarrow t + 1$ and go to 2.

It has been proven that the solution given by the above SA algorithm converges in probability to the global optimal solution of a given problem, under the following conditions [61]:

4. Statistical test on the estimated bicluster structure of a relational data matrix

- (a) **Irreducibility:** we call that the state y is *reachable* at height E from state x if $x = y$ or a sequence of states $x = x_1, x_2, \dots, x_p = y$ exists such that (1) $R(x_t, x_{t+1}) > 0$, for all $t \in \{1, \dots, p-1\}$, and (2) $f(x_t) \leq E$, for all $t \in \{1, \dots, p\}$. We simply call that y is reachable from x if y is reachable from x at some height E . The first condition is that for any pair of states (x, y) , y is reachable from x .
- (b) **Weak reversibility:** The second condition is that, for any $E \in \mathbb{R}$ and for any pair of states (x, y) , y is reachable at height E from x iff x is reachable at height E from y .
- (c) We call that state x is a *local minimum* if no state $y \in \mathcal{S}$ satisfying $f(y) < f(x)$ (i.e., a better solution) is reachable at height $f(x)$ from x . In other words, to find a better solution from a local minimum x , we need to pass through some “worse” states, where the value of the objective function is larger than that of x . We define that the *depth* of a local minimum x is $+\infty$ if x is a global optimal state; otherwise, it is the minimum $E > 0$ such that some state y (i.e., better solution) with $f(y) < f(x)$ exists and y is reachable at height $f(x) + E$ from x . The third condition is that the cooling schedule of temperature satisfies the following conditions: (1) $T_t \geq T_{t+1}$, for all $t \geq 0$, (2) $\lim_{t \rightarrow \infty} T_t = 0$, and (3) $\sum_{t=0}^{\infty} \exp\left(-\frac{d^*}{T_t}\right) = +\infty$, where d^* is the maximum depth of all the states that are locally, but not globally, optimal solutions.

Algorithm 1 is the SA algorithm for obtaining an approximated solution for the optimal set of cluster memberships \hat{g} in terms of the squared residue. In this algorithm, from (4.7), we define that the set of states \mathcal{S} and the objective function f are given by $\mathcal{S} \equiv \mathcal{G}_{K_0 H_0}$ and $f(g) \equiv \mathbf{x}^\top E^{(g)} \mathbf{x}$, respectively. In each step of the algorithm, neighbors $N(g)$ of the current state g are defined as a set of all the cluster memberships that differ from g in exactly one row or column. It must be noted that the size of such neighbors is $|N(g)| = n(K_0 - 1) + p(H_0 - 1)$. We choose a neighbor g' from the uniform distribution on $N(g)$ (i.e., with probability $R(g, g') = 1/|N(g)|$). By these definitions, Algorithm 1 satisfies the conditions of (a) irreducibility and (b) weak reversibility.

Algorithm 2 is the SA algorithm used for finding an approximated solution of the cluster memberships \tilde{g} , which determines the truncation interval. In this algorithm, we define that the set of states \mathcal{S} and the objective function f are given by $\mathcal{S} \equiv \mathcal{G}_{K_0 H_0}$ and $f(g) \equiv \frac{-b^{(g, \hat{g})} - \sqrt{(b^{(g, \hat{g})})^2 - 4a^{(g, \hat{g})}c^{(g, \hat{g})}}{2a^{(g, \hat{g})}}$, respectively. Unlike Algorithm 1, we have to consider the *feasibility* of a solution \tilde{g} , that is, it should satisfy $a^{(\tilde{g}, \hat{g})} < 0$. To guarantee the condition of (a) irreducibility while avoiding infeasible solutions, we defined the neighbors $N(g)$ of the current state g as $N(g) \equiv \mathcal{G}_{K_0 H_0}$. By this definition, for any pair of states (g, g') , transition from g to g' is possible with non-zero probability: $R(g, g') > 0$. Accordingly, we restrict the significant change in the state by controlling the transition probability $R(g, g')$, as in (4.33). By setting the objective function values for infeasible solutions at $+\infty$, we can avoid moving to them throughout the algorithm while satisfying the conditions of (a) irreducibility and (b) weak reversibility.

4. Statistical test on the estimated bicluster structure of a relational data matrix

Algorithm 1 SA algorithm for finding the minimum squared residue solution \hat{g} .

Require: A cooling schedule of temperature $\{T_t\}_{t=0}^{\infty}$ and a threshold ϵ .

Ensure: Approximated optimal set of cluster memberships \hat{g} in terms of the squared residue.

- 1: $t \leftarrow 0$.
 - 2: Randomly generate initial cluster memberships: $\hat{g} = (\hat{g}^{(1)}, \hat{g}^{(2)})$.
 - 3: Compute the initial value of the objective function: $f \leftarrow \mathbf{x}^\top E^{(\hat{g})} \mathbf{x}$.
 - 4: **while** $T_t \geq \epsilon$ **do**
 - 5: Randomly choose a row or column index m from the uniform distribution on $\{1, \dots, n + p\}$.
 - 6: **if** $m \leq n$ **then**
 - 7: $i \leftarrow m$.
 - 8: Randomly generate a new cluster index k' of the i th **row** from the uniform distribution on $\{1, \dots, K\} \setminus \hat{g}_i^{(1)}$. Let \hat{g}' be the set of cluster memberships given by $\hat{g}' = ((\hat{g}')^{(1)}, (\hat{g}')^{(2)})$, $(\hat{g}')_i^{(1)} = k'$, $(\hat{g}')_{i'}^{(1)} = \hat{g}_i^{(1)}$, for $i' \neq i$, and $(\hat{g}')^{(2)} = \hat{g}^{(2)}$.
 - 9: **else**
 - 10: $j \leftarrow m - n$.
 - 11: Randomly generate a new cluster index h' of the j th **column** from the uniform distribution on $\{1, \dots, H\} \setminus \hat{g}_j^{(2)}$. Let \hat{g}' be the set of cluster memberships given by $\hat{g}' = ((\hat{g}')^{(1)}, (\hat{g}')^{(2)})$, $(\hat{g}')_j^{(2)} = h'$, and $(\hat{g}')_{j'}^{(2)} = \hat{g}_j^{(2)}$, for $j' \neq j$.
 - 12: **end if**
 - 13: Compute the value of the objective function: $f' \leftarrow \mathbf{x}^\top E^{(\hat{g}')} \mathbf{x}$.
 - 14: $\Delta f \leftarrow f' - f$.
 - 15: **if** $\Delta f < 0$ **then**
 - 16: $\hat{g} \leftarrow \hat{g}'$, $f \leftarrow f'$.
 - 17: **else**
 - 18: With probability $\exp\left(-\frac{\Delta f}{T_t}\right)$, $\hat{g} \leftarrow \hat{g}'$, $f \leftarrow f'$.
 - 19: **end if**
 - 20: $t \leftarrow t + 1$.
 - 21: **end while**
-

4. Statistical test on the estimated bicluster structure of a relational data matrix

Algorithm 2 SA algorithm for finding the solution \tilde{g} of the truncation interval.

Require: Optimal set of cluster memberships \hat{g} in terms of the squared residue, a cooling schedule of temperature $\{T_t\}_{t=0}^{\infty}$, and a threshold ϵ .

Ensure: Approximated optimal set of cluster memberships \tilde{g} for determining the truncation interval.

- 1: $t \leftarrow 0$.
- 2: Randomly generate initial cluster memberships: $\tilde{g} = (\tilde{g}^{(1)}, \tilde{g}^{(2)})$.
- 3: Compute the initial value of the objective function: if $a^{(\tilde{g}, \hat{g})} = 0$, then $f \leftarrow +\infty$; otherwise, $f \leftarrow (-b^{(\tilde{g}, \hat{g})} - \sqrt{(b^{(\tilde{g}, \hat{g})})^2 - 4a^{(\tilde{g}, \hat{g})}c^{(\tilde{g}, \hat{g})}})/(2a^{(\tilde{g}, \hat{g})})$.
- 4: **while** $T_t \geq \epsilon$ **do**
- 5: Randomly choose the size s of a subset of row or column indices from $\{1, \dots, n + p\}$:

$$\begin{cases} s \leftarrow 1 & \text{with probability } \frac{1}{2} + \frac{1}{2^{n+p}}, \\ s \leftarrow s' & \text{with probability } \frac{1}{2^{s'}}, \text{ for } s' \in \{2, \dots, n + p\}. \end{cases} \quad (4.33)$$

- 6: Randomly choose a set of s row or column indices \mathcal{S} without duplication from the uniform distribution.
 - 7: $\tilde{g}' \leftarrow \tilde{g}$.
 - 8: **for** each row or column index in \mathcal{S} **do**
 - 9: **if** the i th row is selected **then**
 - 10: Randomly generate a new cluster index k' of the i th row from the uniform distribution on $\{1, \dots, K\} \setminus \tilde{g}_i^{(1)}$. $(\tilde{g}')_i^{(1)} \leftarrow k'$.
 - 11: **else if** the j th column is selected **then**
 - 12: Randomly generate a new cluster index h' of the j th column from the uniform distribution on $\{1, \dots, H\} \setminus \tilde{g}_j^{(2)}$. $(\tilde{g}')_j^{(2)} \leftarrow h'$.
 - 13: **end if**
 - 14: **end for**
 - 15: Compute the value of the objective function: if $a^{(\tilde{g}', \hat{g})} = 0$, then $f \leftarrow +\infty$; otherwise, $f' \leftarrow (-b^{(\tilde{g}', \hat{g})} - \sqrt{(b^{(\tilde{g}', \hat{g})})^2 - 4a^{(\tilde{g}', \hat{g})}c^{(\tilde{g}', \hat{g})}})/(2a^{(\tilde{g}', \hat{g})})$.
 - 16: $\Delta f \leftarrow f' - f$.
 - 17: **if** $\Delta f < 0$ **then**
 - 18: $\tilde{g} \leftarrow \tilde{g}', f \leftarrow f'$.
 - 19: **else**
 - 20: With probability $\exp\left(-\frac{\Delta f}{T_t}\right)$, $\tilde{g} \leftarrow \tilde{g}', f \leftarrow f'$.
 - 21: **end if**
 - 22: $t \leftarrow t + 1$.
 - 23: **end while**
-

4. Statistical test on the estimated bicluster structure of a relational data matrix

Regarding the cooling schedule of temperature, we can use the following definition [61], which satisfies the conditions (1), (2), and (3) in (c):

$$T_t = c / \log(t + 2) \quad \text{for all } t \geq 0, \quad (4.34)$$

where c is a constant satisfying $c \geq d^*$. In our cases, for instance, we can define the constant c as follows:

$$c \equiv \|\mathbf{x}\|_2^2 - \frac{1}{np} \left([1 \cdots 1] \mathbf{x} \right)^2, \quad (4.35)$$

for Algorithm 1, since $d^* \leq \max_{g \in \mathcal{G}_{K_0 H_0}} \mathbf{x}^\top E^{(g)} \mathbf{x} - \min_{g \in \mathcal{G}_{K_0 H_0}} \mathbf{x}^\top E^{(g)} \mathbf{x} \leq c$. Here, we take into account the fact that the cluster memberships $\underline{g} \equiv \arg \max_{g \in \mathcal{G}_{K_0 H_0}} \mathbf{x}^\top E^{(g)} \mathbf{x}$ are attained by assigning all the elements of an observed matrix into a single block, where the objective function value is given by $\mathbf{x}^\top E^{(\underline{g})} \mathbf{x} = \|\mathbf{x}\|_2^2 - \frac{1}{np} \left([1 \cdots 1] \mathbf{x} \right)^2$, and that $\min_{g \in \mathcal{G}_{K_0 H_0}} \mathbf{x}^\top E^{(g)} \mathbf{x} \geq 0$. It must be noted that, in Algorithm 2, there is no state that is local but not global minimum (i.e., all local minima are also global minima); this is because, for any pair of states (g, g') , g' is reachable at height $f(g)$ from g . Therefore, from the result in [61], the convergence in probability to a global minimum state is guaranteed without the condition (3) in (c).

Practically, an algorithm based on the cooling schedule (4.34) is too slow, that is, it requires much computation time before convergence. Therefore, in the experiments in Section 4.4, we used the cooling schedule of $T_t = T_0 \times r^t$ with $r < 1$ for all $t \geq 0$, though this definition satisfies only the conditions (1) and (2), not (3).

4.4 Experiments

To show the validity of our proposed test, we compared its behavior with that of a *naive* statistical test, which does not consider the selection event. By ignoring the fact that the set of cluster memberships \hat{g} was selected based on the data vector \mathbf{x} , we construct a naive test (**which is invalid in fact**) with test statistic T in (4.13) by assuming

$$T | \{\mathbf{z}, \mathbf{u}\} \sim \chi_{(np - K_0 H_0)}, \quad (4.36)$$

from (4.21). The p -value of such a naive test is given by

$$p_T = \begin{cases} 1 - \frac{\gamma\left(\frac{np - K_0 H_0}{2}, \frac{T^2}{2}\right)}{\Gamma\left(\frac{np - K_0 H_0}{2}\right)} & \text{if } 0 \leq T, \\ 0 & \text{otherwise,} \end{cases} \quad (4.37)$$

where $\Gamma(\cdot)$ is the Gamma function.

4. Statistical test on the estimated bicluster structure of a relational data matrix

In the following sections 4.4.1 and 4.4.2, we check the behavior of the p -values and the TPR and FPR when using the proposed and naive tests, in order to show the validity of our proposed method. In these sections, we use the term “realizable case” to indicate that the hypothetical cluster numbers of rows and columns (K_0, H_0) are equal to the null ones (K, H) , and the term “unrealizable case” to indicate that at least one of $K_0 < K$ and $H_0 < H$ holds, as described in Section 4.1.

Aside from the experiments in this section, we conducted sensitivity analysis of the approximated version of the proposed test with respect to the cooling schedule of SA in the realizable case in Appendix 4.E. Moreover, we conducted an additional experiment to employ an existing fast biclustering method [140] instead of the proposed SA-based algorithm for estimating the cluster memberships in Appendix 4.F.

4.4.1 Exact test in realizable case: $(K_0, H_0) = (K, H)$

First, we check the behavior of the p -values calculated by using the proposed (4.31) and naive (4.37) tests, under the condition that the given set of cluster numbers (K_0, H_0) are equal to that of the null one (K, H) . As shown in Section 4.3.2, the p -value of the proposed test follows the uniform distribution on $[0, 1]$, while there is no such guarantee for that of the naive test.

For experiment, we independently generated data matrices with the sizes of $(n, p) = (5, 5), (6, 6), \dots, (9, 9)$. We set the null and hypothetical sets of cluster numbers at $(2, 2)$; we defined the null cluster memberships as $g_i^{(N),(1)} = (i \bmod K) + 1$, for all i , and $g_j^{(N),(2)} = (j \bmod H) + 1$, for all j . In regard to the mean vector $\boldsymbol{\mu}_0$, we tried the following five settings:

$$\boldsymbol{\mu}_0^{(l)} = \left(1 - \frac{l-1}{5}\right) \left[\text{vec} \left(\begin{bmatrix} 0.7 & 0.55 \\ 0.5 & 0.6 \end{bmatrix} \right) - 0.5 \right] + 0.5, \quad l = 1, \dots, 5. \quad (4.38)$$

Based on the above settings, we generated 1000 data vectors by $\boldsymbol{x} \sim N(\boldsymbol{\mu}_0^{(l)}, 0.05^2 I)$, for each setting of matrix size (n, p) and mean vector $\boldsymbol{\mu}_0$. Figure 2 shows the examples of the generated data matrices. For each generated data vector \boldsymbol{x} , we computed the squared residues of all the patterns of cluster memberships g . Subsequently, we chose the optimal set of cluster memberships \hat{g} (i.e., solution with the minimum squared residue) and checked if it is equivalent to the null set of cluster memberships $g^{(N)}$. For both cases of $\hat{g} = g^{(N)}$ and $\hat{g} \neq g^{(N)}$, we computed the test statistic T in (4.13), the truncated interval in (4.28), and the p -values in (4.31) and (4.37). Subsequently, we plotted the results as follows:

- For the trials where $\hat{g} = g^{(N)}$ holds (i.e., under null hypothesis), we plotted the p -values given by (4.31) and (4.37), in Figures 3 and 4, respectively. We also plotted (i) the test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed and naive tests and (ii) the accuracy of the clustering algorithm \mathcal{A} , that

4. Statistical test on the estimated bicluster structure of a relational data matrix

is, the ratio of the number of such null cases (i.e., $\hat{g} = g^{(N)}$) to the 1000 trials, for each setting, in Figures 5 and 6, respectively.

- For null (i.e., $\hat{g} = g^{(N)}$) and alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, we plotted the FPR and TPR, in Figure 7, respectively. We also plotted the AUC score in Figure 8.

From Figures 3, 4, and 5, we see that the distribution of the p -values of the proposed test was closer to the uniform distribution on $[0, 1]$ than that of the naive test, particularly when the difference in block-wise mean between the blocks was small. This result shows that the proposed test can successfully take into account the selective bias of using the squared residue minimization solution, by using the truncated chi distribution in (4.13). However, the p -values of the naive test based on (4.36) did not follow the uniform distribution on $[0, 1]$, since we did not treat the selective bias and conducted tests based on the (not truncated) chi distribution in the naive test. It must be noted that, in our problem setting, unlike the common statistical tests, the assertion of the null hypothesis ($(g^{(N),(1)}, g^{(N),(2)}) = (\hat{g}^{(1)}, \hat{g}^{(2)})$) is stronger than that of alternative hypothesis ($(g^{(N),(1)}, g^{(N),(2)}) \neq (\hat{g}^{(1)}, \hat{g}^{(2)})$). This results in that the p -values of the naive test are biased toward **larger** values than the correct ones.

In regard to the test performance, from the results in the top of Figure 7, we see that the FPR was low in all the settings (i.e., proposed and naive; significance rate $\alpha = 0.1, 0.05$, and 0.01 ; and block-wise mean μ_0). The results in the bottom of Figure 7 shows that the TPR of the proposed test was higher than that of the naive test in the same setting, which is consistent with the discussion in the previous paragraph. However, the TPR of the proposed test was not sufficiently close to one, in all the settings. This can be attributed to the few “true positive” cases in the realizable setting. Figure 6 shows that the estimated block structure \hat{g} that attained the minimum squared residue was equivalent to the null one $g^{(N)}$ in most cases. In other words, almost all trials were “null cases,” where the small number of false negative cases significantly affect the TPR. Particularly, when there is an increase in the matrix size or in the difference in the block-wise mean between the blocks, it becomes easier to estimate the null block structure, and the clustering algorithm almost always outputs the correct cluster memberships. With regard to the AUC score, from Figure 8, we see that the proposed test outperformed the naive one in all the settings.

4.4.2 Exact test in unrealizable cases: $K_0 < K$ or $H_0 < H$

Next, we compared the behavior of the proposed and naive tests in the unrealizable cases, that is, either $K_0 < K$ or $H_0 < H$ holds.

For the experiment, we independently generated data matrices with the sizes of $(n, p) = (5, 5), (6, 6), \dots, (9, 9)$. We set the null set of cluster numbers at $(3, 2)$ and defined the null cluster memberships as in Section 4.4.1. In regard to the mean vector μ_0 , we tried the

4. Statistical test on the estimated bicluster structure of a relational data matrix

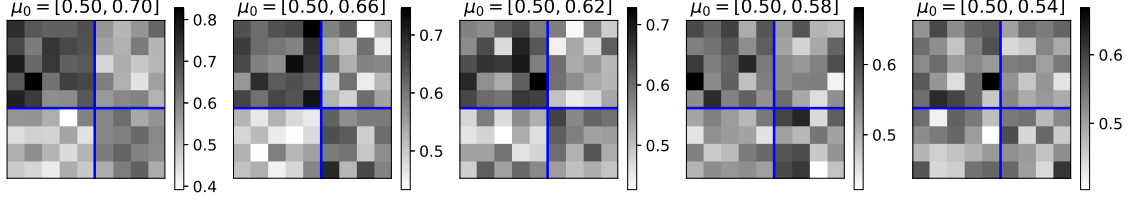


Figure 2: Examples of the observed data matrices with the size of $(n, p) = (9, 9)$, which are generated based on the different block-wise means. The title of each figure shows the range of the block-wise mean vector μ_0 . The blue lines show the **null** cluster memberships. For visibility, we plotted the matrices whose rows and columns were sorted according to their null clusters.

following five settings:

$$\mu_0^{(l)} = \left(1 - \frac{l-1}{5}\right) \left[\text{vec} \left(\begin{bmatrix} 0.7 & 0.55 \\ 0.5 & 0.6 \\ 0.55 & 0.5 \end{bmatrix} \right) - 0.5 \right] + 0.5, \quad l = 1, \dots, 5. \quad (4.39)$$

Based on the above settings, we generated 1000 data vectors by $\mathbf{x} \sim N(\mu_0^{(l)}, 0.05^2 I)$, for each setting of matrix size (n, p) and mean vector μ_0 . For each generated data vector \mathbf{x} , we computed the squared residues of all the patterns of cluster memberships g . Subsequently, we chose the optimal set of cluster memberships \hat{g} with a given set of cluster numbers (K_0, H_0) . In regard to the hypothetical cluster numbers, we tried the following five settings: $(K_0, H_0) = (1, 1), (2, 1), (3, 1), (1, 2)$, and $(2, 2)$. For each setting, based on the selected result \hat{g} , we computed the test statistic T in (4.13), the truncated interval in (4.28), and the p -values in (4.31) and (4.37). Finally, we plotted the TPR of the proposed and naive tests in Figure 9.

Figure 9 shows that the TPR of the proposed test was higher than that of the naive test in the same setting; however, in most cases, there was a small difference between them. This may be attributed to the fact that we set the matrix size (n, p) and the hypothetical block size (K_0, H_0) at small numbers in order to perform the exact test, which is computationally expensive, and thus there is a marginal effect of selecting the optimal block structure \hat{g} from all the patterns $\mathcal{G}_{K_0 H_0}$. It must be noted that, unlike the realizable case in Section 4.4.1, the block structures output by the clustering algorithm were **always** different from the null ones because the hypothetical set of cluster numbers were insufficient to represent the null block structure in the unrealizable cases. In other words, all the 1000 trials in each setting correspond to the alternative cases. The TPRs of the proposed and naive tests were comparable particularly under the following two settings: (1) the case where we set the hypothetical number of blocks at $(K_0, H_0) = (1, 1)$ and (2) the case where the difference in the null block-wise mean μ between the blocks was relatively big. These results were caused by the nature of the biclustering problem itself, as well as by the limitation in the

4. Statistical test on the estimated bicluster structure of a relational data matrix

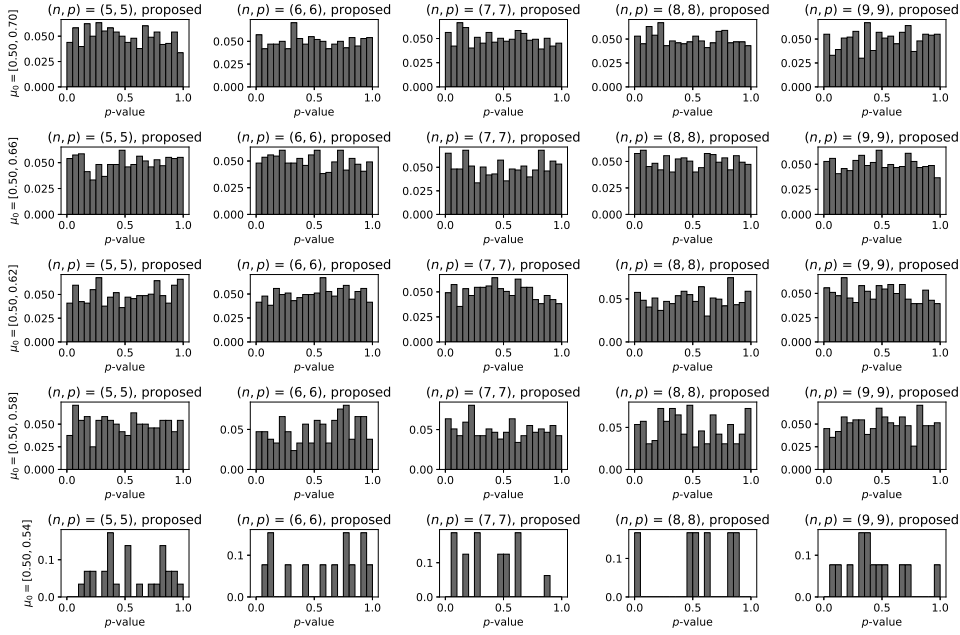


Figure 3: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **proposed** test (4.31) on the set of cluster memberships \hat{g} with the minimum squared residue.

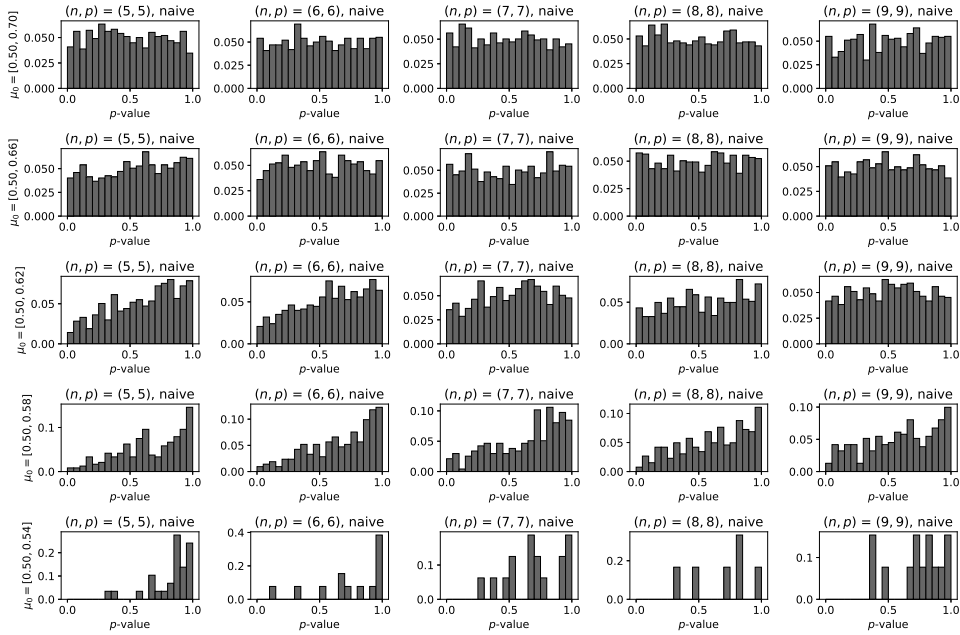


Figure 4: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **naive** test (4.37).

4. Statistical test on the estimated bicluster structure of a relational data matrix

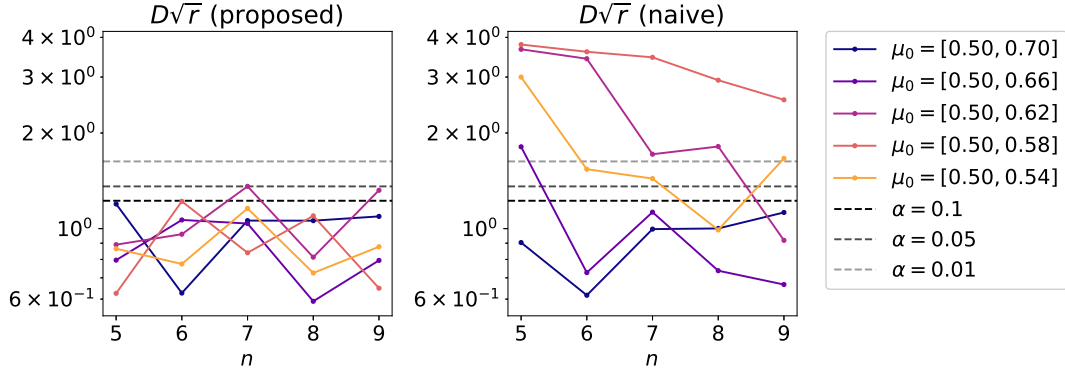


Figure 5: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) tests. The null hypothesis that p -value follows the uniform distribution on $[0, 1]$ is rejected if $D\sqrt{r} > \alpha$, where α is a given significance level.

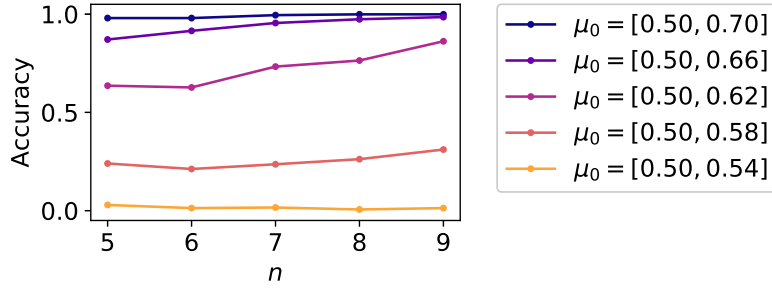


Figure 6: The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 . For this experiment, we used the setting of $n = p$.

power of the proposed selective test. In the case of (1), we assume that the entire data matrix A consists of a single block, and thus there is only one possible estimated bicluster structure \hat{g} , regardless of the selection event. Therefore, in this case, the proposed and naive tests are equivalent in the first place. As for the case of (2), since the difference in the null block-wise mean between the blocks was big, the test statistics of both the proposed and naive tests got large enough for the null hypothesis to be rejected.

4.4.3 Approximated test in both realizable and unrealizable cases

Finally, we checked the behavior of the approximated test introduced in Section 4.3.3. In both realizable and unrealizable cases, we generated data matrices with the sizes of $(n, p) = (10 + 2 \times m, 10 + 2 \times m)$, for $m = 0, 1, \dots, 4$, in the same way as that in Sections 4.4.1 and 4.4.2. Concerning the following conditions, we used the same setting as in Sections 4.4.1 and 4.4.2, respectively, for the realizable and unrealizable cases: the null and hypothetical sets of cluster numbers, the definition of null cluster memberships $g^{(N)}$,

4. Statistical test on the estimated bicluster structure of a relational data matrix

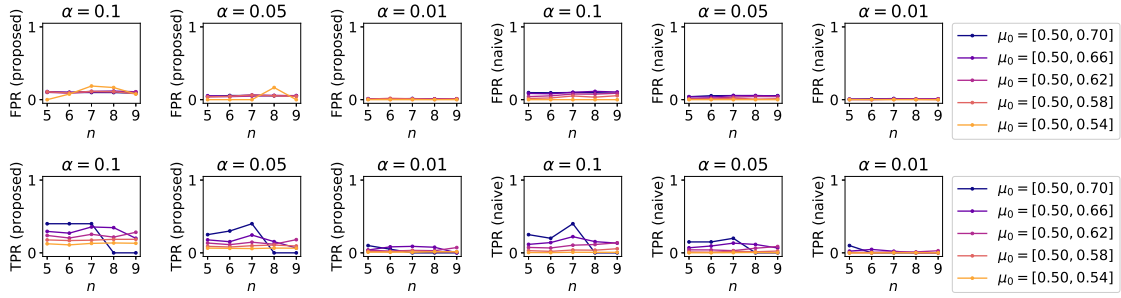


Figure 7: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the proposed (left) and naive (right) statistical tests. If there were no null (i.e., $\hat{g} = g^{(N)}$) or alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, respectively, then the corresponding points of FPR or TPR would not have been plotted.

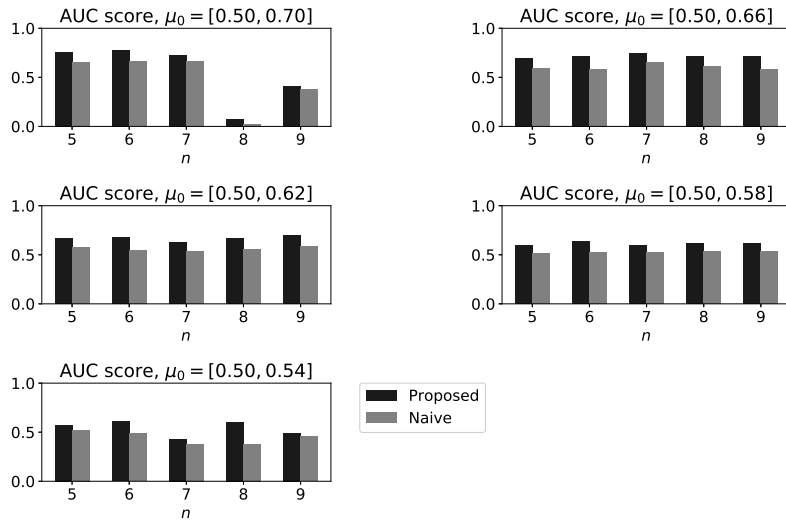


Figure 8: AUC score in the realizable case for the proposed and naive statistical tests.

4. Statistical test on the estimated bicluster structure of a relational data matrix

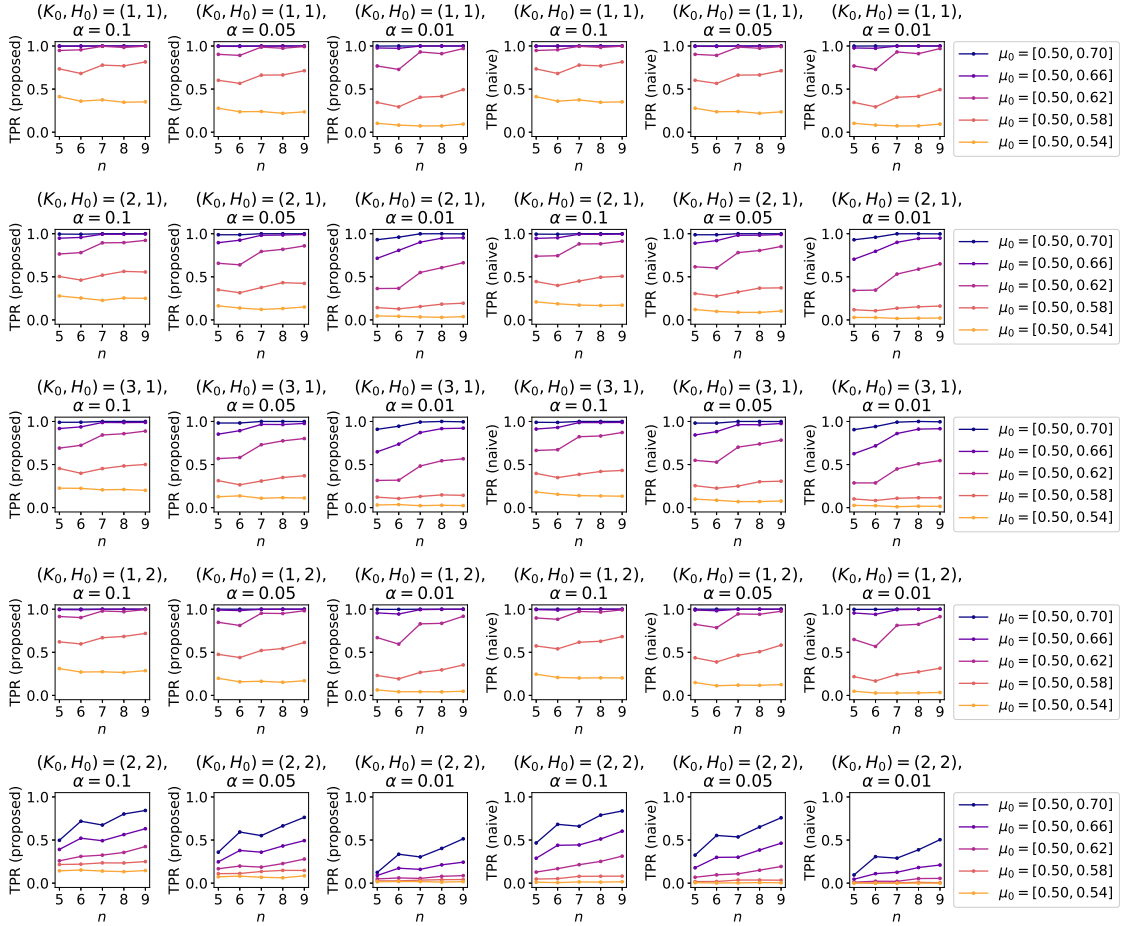


Figure 9: TPR in the unrealizable case (i.e., $K_0 < K$ or $H_0 < H$) with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the proposed (left) and naive (right) statistical tests.

4. Statistical test on the estimated bicluster structure of a relational data matrix

mean vectors and the standard deviation σ_0 , and the number of data vectors for each setting. Concerning the SA algorithms, both in the Algorithms 1 and 2, we defined the cooling schedule as follows: $T_0 = 10$, $T_t = T_0 \times 0.99^t$ for all t .

As in the cases of the exact tests, Figures 10 and 11, respectively, show the histograms of the p -values of the proposed and naive approximated tests in the realizable case. For the realizable case, we also plotted (i) the test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36], for the p -values of the proposed and naive tests, and (ii) the accuracy of the approximated clustering algorithm in Figures 12 and 13, respectively. Figure 14 shows the FPR and TPR in the realizable case, Figure 15 shows the AUC score in the realizable case, and Figure 16 shows the TPR in the unrealizable cases.

Figures 10, 11, and 12 show that the distributions of the p -values of the proposed test were closer to the uniform distribution on $[0, 1]$ than that of the naive test, as in the result of the exact test in Section 4.4.1. Concerning the test performance in the realizable case, Figure 14 shows that the FPR was low in all the settings, and the TPR of the proposed test was higher than that of the naive test in the same setting. However, as in the exact case, the TPR of the proposed test was not sufficiently close to one in all the setting; this can be attributed to the few “true positive” cases. In the next Section 4.4.4, we checked more difficult cases for the approximated clustering algorithm, where the null cluster numbers were more than $(2, 2)$. With regard to the AUC score, from Figure 15, we see that the proposed test achieved comparable or better performance than the naive one in all the settings.

4.4.4 Approximated test in the realizable case, $(K, H) = (3, 3), (4, 4), (5, 5)$

To check the behavior of the p -values, FPR, and TPR of the proposed test in more difficult settings, where the clustering algorithm cannot successfully estimate the cluster memberships in most cases, we tried the following three settings of null cluster numbers: $(K, H) = (3, 3), (4, 4)$, and $(5, 5)$. These settings have more patterns of the possible block structures than those in the case of $(K, H) = (2, 2)$ in Section 4.4.3. Hence, it becomes difficult for the approximated clustering algorithm (which stops at a fixed finite number of steps in the experiment) to output the null set of cluster memberships.

We generated data matrices in the same way as that in Section 4.4.3. Concerning the following conditions, we used the same setting as that of the realizable case in Section 4.4.3: the set of matrix sizes (n, p) , the definition of the null cluster memberships $g^{(N)}$, the standard deviation σ_0 , and the cooling schedule of the SA algorithm. We tried the following three settings of the null number of blocks: $(3, 3), (4, 4)$, and $(5, 5)$; subsequently, for each

4. Statistical test on the estimated bicluster structure of a relational data matrix

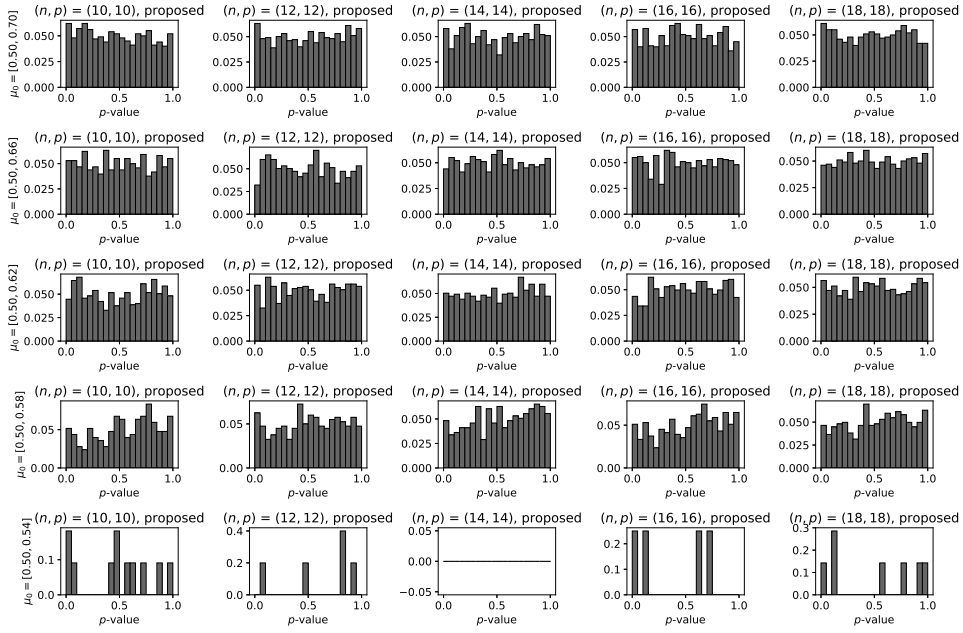


Figure 10: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **approximated** version of the **proposed** test.

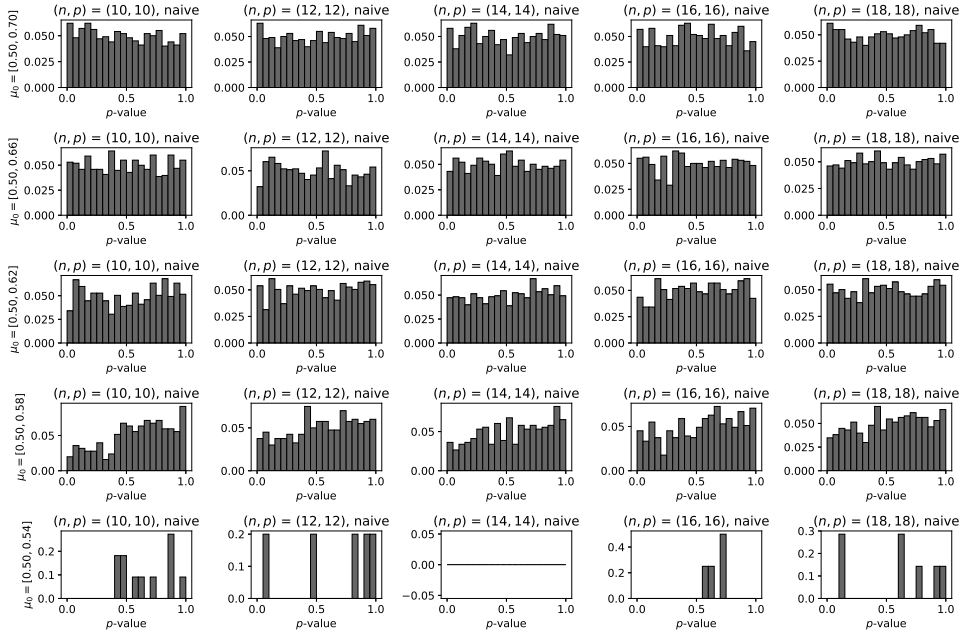


Figure 11: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **approximated** version of the **naive** test (4.37).

4. Statistical test on the estimated bicluster structure of a relational data matrix

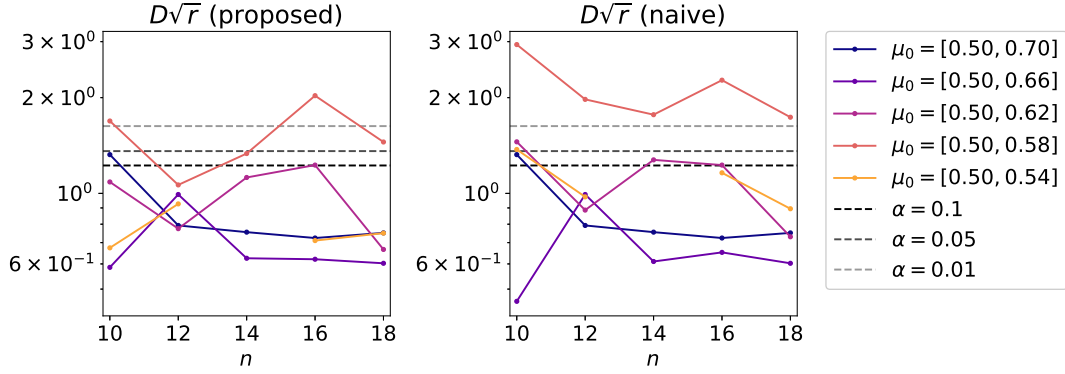


Figure 12: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) **approximated** tests.

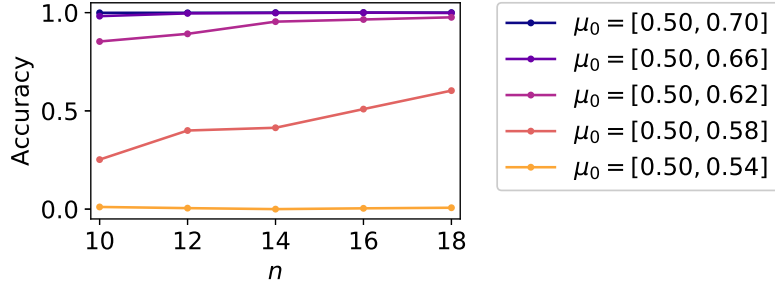


Figure 13: The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the **approximated** clustering algorithm in Section 4.3.3. For the experiment, we used the setting of $n = p$.

setting, we defined the mean vector μ_0 as follows:

$$\mu_0^{(l)} = \left(1 - \frac{l-1}{5}\right) \left[\text{vec} \left(\begin{bmatrix} 0.6 & 0.55 & 0.7 \\ 0.4 & 0.6 & 0.5 \\ 0.65 & 0.5 & 0.6 \end{bmatrix} \right) - 0.5 \right] + 0.5, \quad l = 1, \dots, 5, \quad (4.40)$$

for $(K, H) = (3, 3)$,

$$\mu_0^{(l)} = \left(1 - \frac{l-1}{5}\right) \left[\text{vec} \left(\begin{bmatrix} 0.6 & 0.55 & 0.7 & 0.5 \\ 0.4 & 0.6 & 0.5 & 0.7 \\ 0.65 & 0.5 & 0.6 & 0.4 \\ 0.5 & 0.4 & 0.45 & 0.6 \end{bmatrix} \right) - 0.5 \right] + 0.5, \quad l = 1, \dots, 5, \quad (4.41)$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

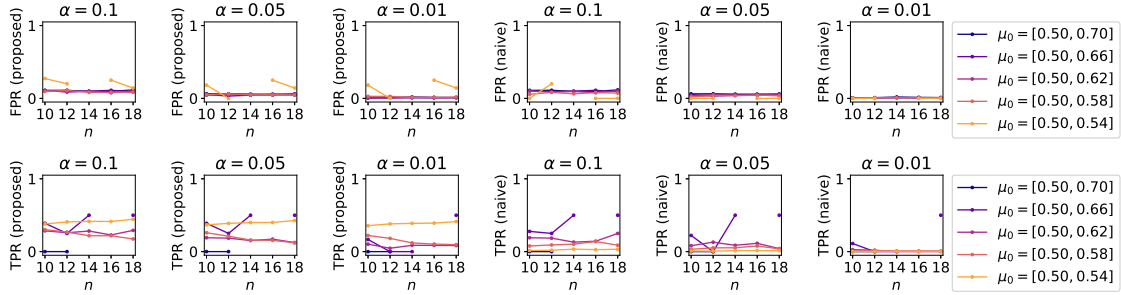


Figure 14: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests.

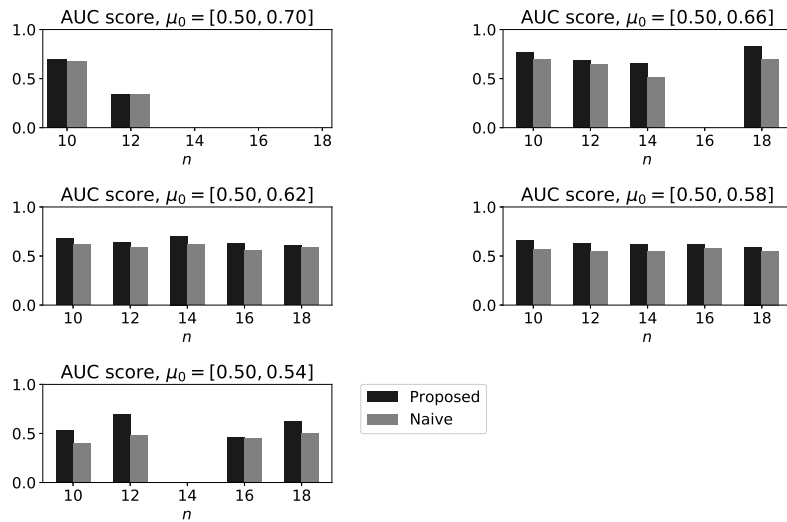


Figure 15: AUC score in the realizable case for the **approximated** version of the proposed and naive statistical tests. If there were no null (i.e., $\hat{g} = g^{(N)}$) or alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, respectively, then the corresponding bars would not have been plotted.

4. Statistical test on the estimated bicluster structure of a relational data matrix

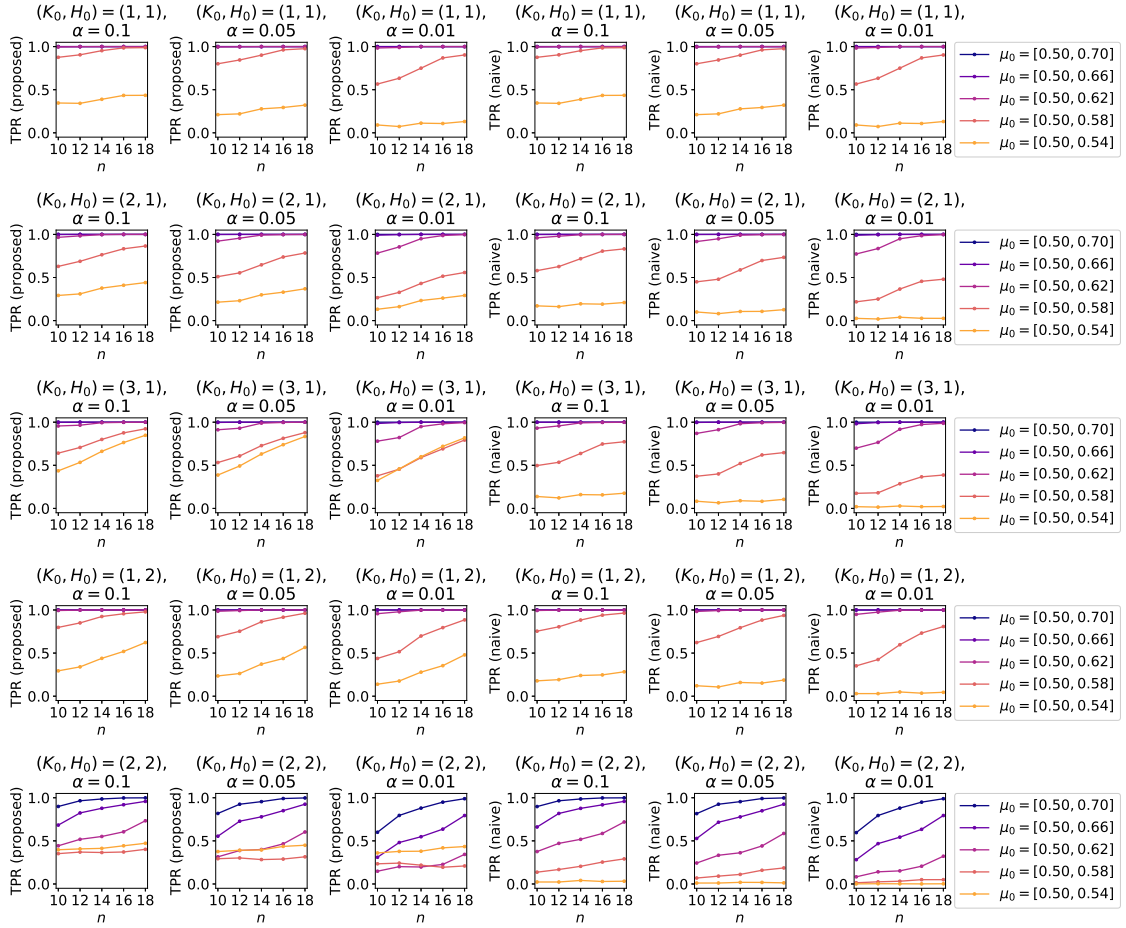


Figure 16: TPR in the unrealizable case (i.e., $K_0 < K$ or $H_0 < H$) with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests.

4. Statistical test on the estimated bicluster structure of a relational data matrix

for $(K, H) = (4, 4)$, and

$$\mu_0^{(l)} = \left(1 - \frac{l-1}{5}\right) \left[\text{vec} \left(\begin{bmatrix} 0.6 & 0.55 & 0.7 & 0.5 & 0.65 \\ 0.4 & 0.6 & 0.5 & 0.7 & 0.55 \\ 0.65 & 0.5 & 0.6 & 0.4 & 0.45 \\ 0.5 & 0.4 & 0.45 & 0.6 & 0.7 \\ 0.7 & 0.65 & 0.55 & 0.45 & 0.6 \end{bmatrix} - 0.5 \right) + 0.5, \right.$$

$$l = 1, \dots, 5, \tag{4.42}$$

for $(K, H) = (5, 5)$. We set the number of data vectors for each setting at 500.

We plotted the test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed and naive tests and the accuracy of the approximated clustering algorithm in Figures 17 and 18, respectively. Figures 19, 20, and 21 show the FPR and TPR of the proposed and naive tests. These figures show that the TPRs of both the proposed and naive tests were higher than the case of $(K, H) = (2, 2)$; the TPR of the proposed test was higher than that of the naive one in these settings. Figures 22, 23, and 24 show the AUC score of the proposed and naive tests. From these results, we see that the proposed test achieved higher AUC score than the naive one in most settings.

4.5 Discussions

In this section, we discuss the following three points about the proposed statistical test: its power, the trade-off between computational efficiency and accuracy, and the extension of the finding to more generalized cases.

First, as also pointed out in a study [103], the null distribution of the test statistic of the proposed test is given by the conditioning on \mathbf{z} and \mathbf{u} , besides the selected set of cluster memberships \hat{g} , which leads to a reduction in the test power [46]. For now, we do not have any way of removing these unnecessary parameters, owing to the problem setting of an LBM. In a one-way clustering problem, where there are n data vectors with p dimensions, we can at least approximate the distributions of \mathbf{z} and \mathbf{u} based on their histograms; however, in the LBM setting, there is only a single observed matrix with the size of $n \times p$. Solving this problem is beyond the scope of this dissertation; future studies should focus on constructing a more powerful selective test on a bicluster structure by using an additional technique such as a bootstrap method [142, 144].

Second, we have proposed both exact and approximated tests to cope with the combinatorial explosion of the possible block memberships. The null distribution (4.13) of the proposed test statistic is based on the assumption that the estimated cluster memberships \hat{g} is the global minimum solution of the squared residue, which is difficult to obtain in the first place. Although it is guaranteed that the solutions of the two SA algorithms 1 and 2 converge in probability to the globally optimal solutions of their corresponding problems,

4. Statistical test on the estimated bicluster structure of a relational data matrix

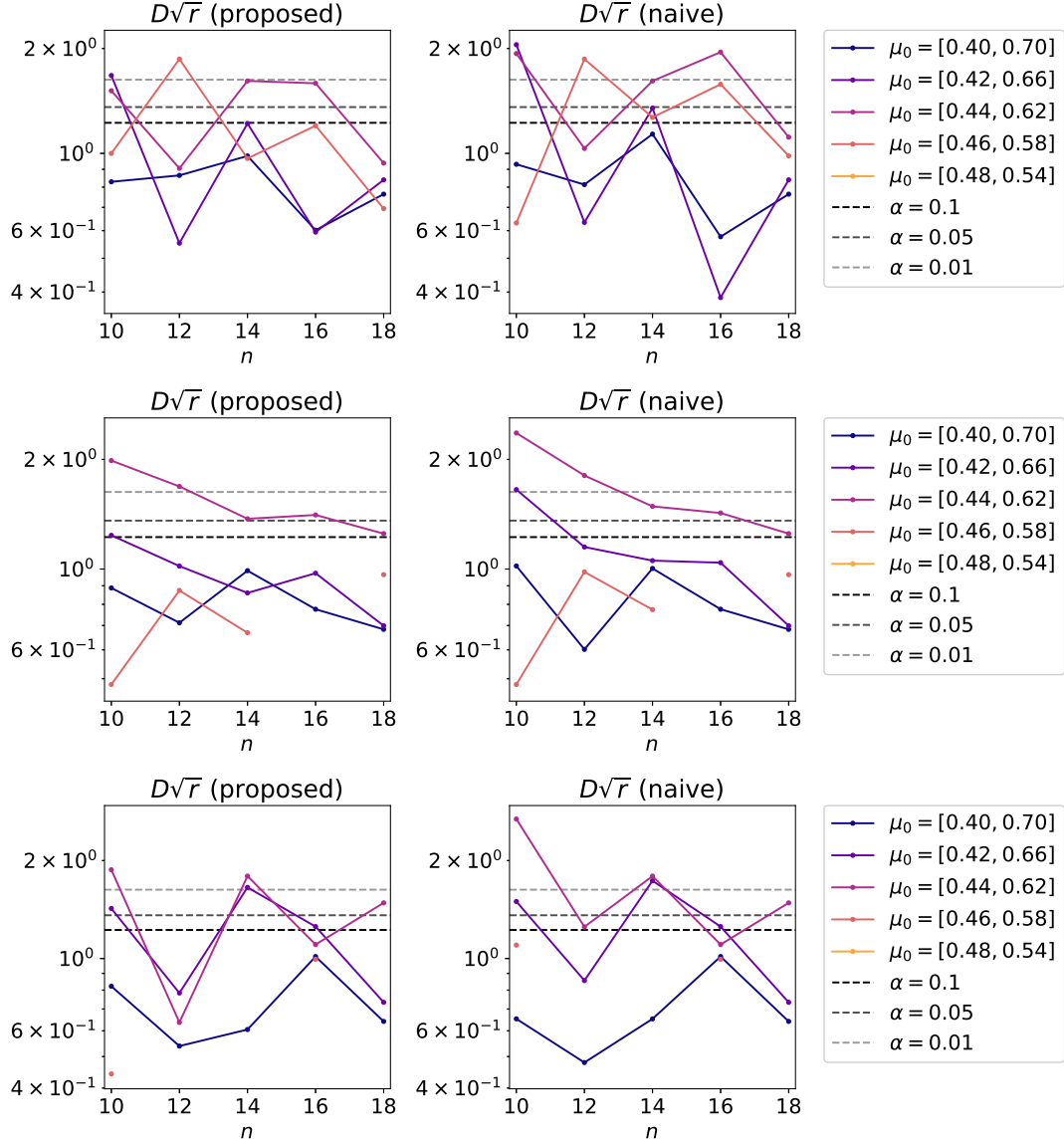


Figure 17: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) **approximated** tests, where $(K, H) = (3, 3)$ (top), $(4, 4)$ (middle), and $(5, 5)$ (bottom).

4. Statistical test on the estimated bicluster structure of a relational data matrix

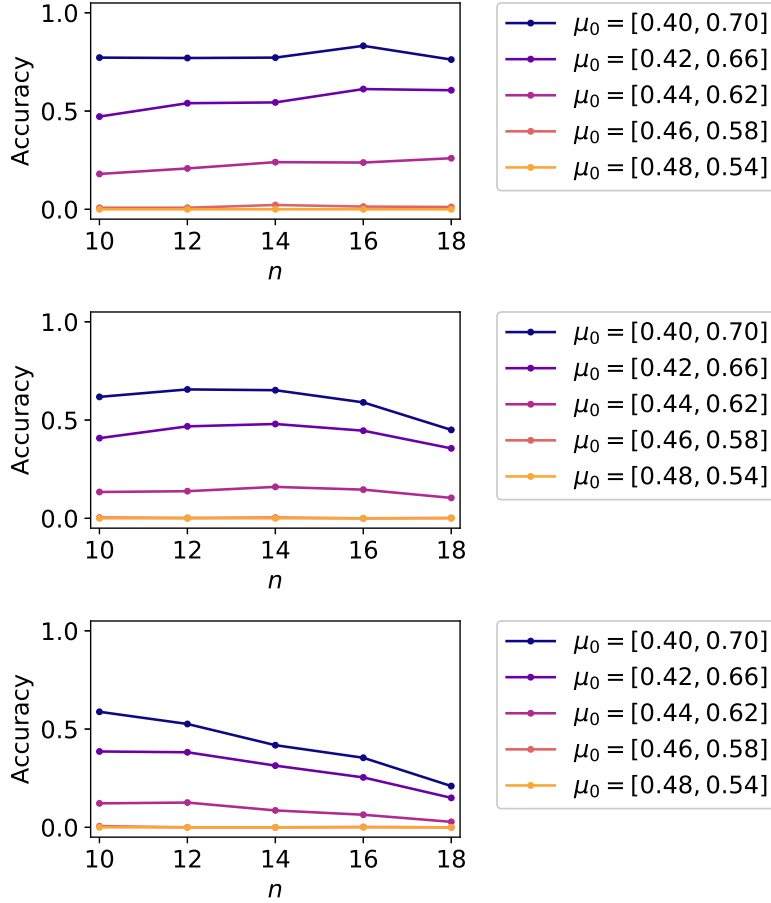


Figure 18: The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$), for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the **approximated** clustering algorithm in Section 4.3.3; $(K, H) = (3, 3)$ (top), $(4, 4)$ (middle), and $(5, 5)$ (bottom). For experiment, we used the setting of $n = p$.

4. Statistical test on the estimated bicluster structure of a relational data matrix

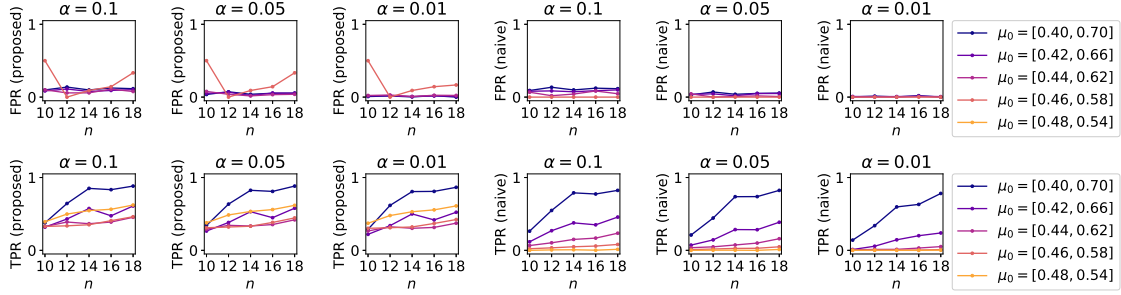


Figure 19: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests, where $(K, H) = (3, 3)$.

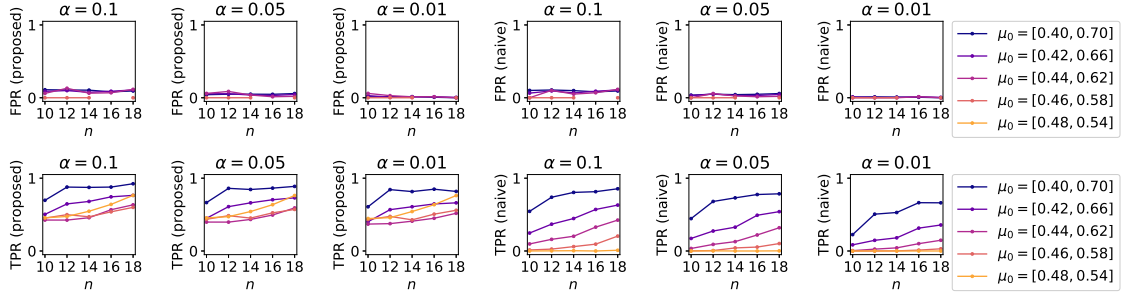


Figure 20: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests, where $(K, H) = (4, 4)$.

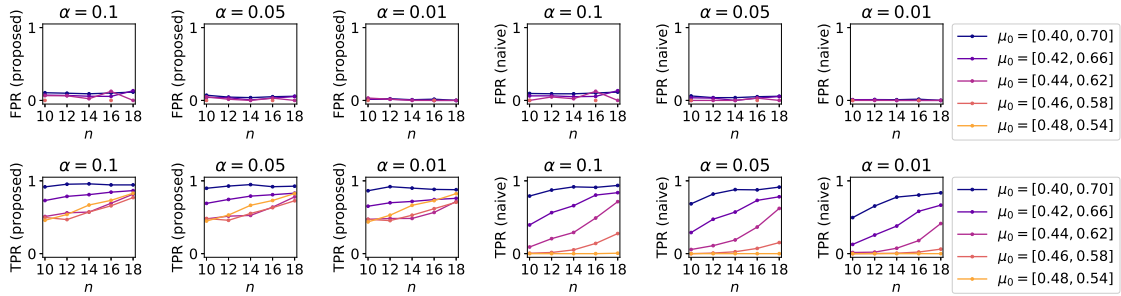


Figure 21: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests, where $(K, H) = (5, 5)$.

4. Statistical test on the estimated bicluster structure of a relational data matrix

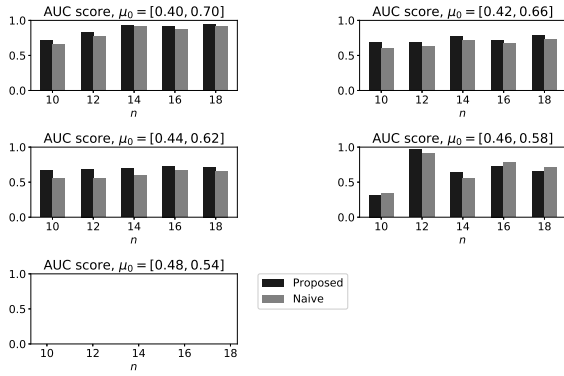


Figure 22: AUC score in the realizable case for the **approximated** version of the proposed and naive statistical tests, where $(K, H) = (3, 3)$.

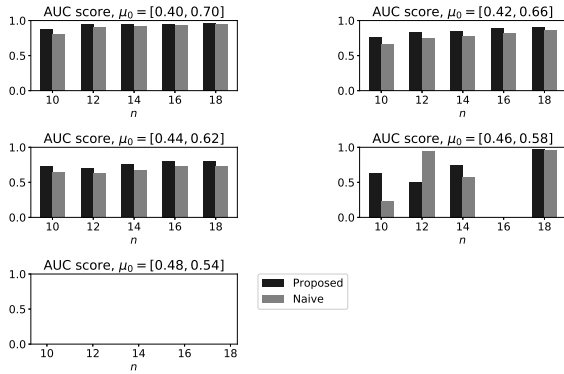


Figure 23: AUC score in the realizable case for the **approximated** version of the proposed and naive statistical tests, where $(K, H) = (4, 4)$.

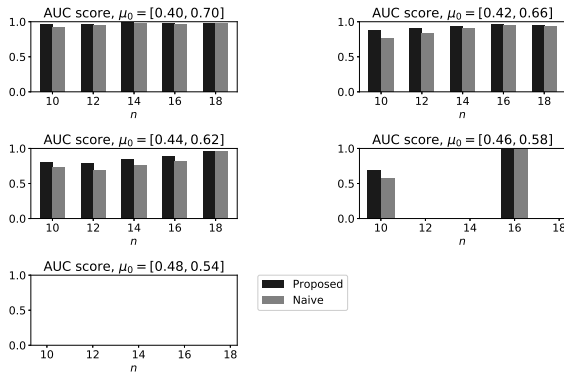


Figure 24: AUC score in the realizable case for the **approximated** version of the proposed and naive statistical tests, where $(K, H) = (5, 5)$.

4. Statistical test on the estimated bicluster structure of a relational data matrix

we cannot validate the outputs of these algorithms with a finite number of steps, which are used in practice. It would be more desirable to derive an *exact* p -value of some other *approximated* test. Stopping the SA algorithms in a constant number of steps would also affect the accuracy of the test; to find the optimal solutions, we should have checked all the patterns of possible block memberships, which increase with the observed matrix size and the number of blocks. However, if we increase the number of steps according to such a problem size, then computation of the SA algorithms will get intractable. Therefore, it would be another important direction to seek a more computationally efficient test, which mitigates this trade-off.

Finally, the proposed test enables us to perform a valid statistical inference for a Gaussian LBM, where we assume that each element of an observed matrix independently follows a Gaussian distribution, given a block structure. This Gaussian assumption is crucial for deriving the exact p -value in the selective inference framework, as in [103]. However, in many practical datasets, including the “MovieLens” dataset of movie ratings [62] and the dataset of document-word relationships in NeurIPS conference papers [120], the elements of the observed matrix take discrete values, where the proposed test cannot be employed. So far, there has been no selective test that can be directly applied to binary data vectors from a Bernoulli distribution. To address this problem, a randomized model selection method [143] has been proposed to construct an asymptotically valid selective test on binary data by adding a random noise to the statistic used for a selection event. By using such a technique, future studies should generalize the proposed test for non-Gaussian cases.

4.6 Chapter conclusion

We developed a new selective inference method on the row and column cluster memberships of a latent block model given by a clustering algorithm based on the squared residue minimization. By considering the selective bias, which is caused by the fact that the hypothetical block structure is estimated based on a given data matrix, we constructed a valid test based on a truncated chi distribution. Since such an exact test required us to obtain the global optimal solutions of two combinatorial optimization problems, we also constructed an approximated test based on simulated annealing algorithms. Experimental results showed that the proposed exact and approximated tests worked successfully, compared to the naive tests that did not take the selective bias into account.

4. Statistical test on the estimated bicluster structure of a relational data matrix

4.A Proof of (4.10) that $E^{(g)} - E^{(g')} \neq O$ for $g, g' \in \mathcal{G}_{K_0 H_0}$, $g \neq g'$

Proof. We prove that $E^{(g)} - E^{(g')} \neq O$ by contradiction. Let $g, g' \in \mathcal{G}_{K_0 H_0}$ be two sets of cluster memberships, both of which have $K_0 \times H_0$ blocks or less and which satisfy $g \neq g'$. Specifically, we denote the exact number of blocks of g as $(K_0^{(g)}, H_0^{(g)})$. Assume that $E^{(g)} - E^{(g')} = O$ holds. Then, for all $\mathbf{x} \in \mathbb{R}^{np}$, we have $\mathbf{x}^\top E^{(g)} \mathbf{x} = \mathbf{x}^\top E^{(g')} \mathbf{x}$. In other words, from (4.6), block structures g and g' yield the same squared residue σ^2 for any data matrix A .

Let us consider a data matrix A that has a block structure g , and all of the elements in the (k, h) th block are $(k-1)H_0^{(g)} + h$, where $k = 1, \dots, K_0^{(g)}$ and $h = 1, \dots, H_0^{(g)}$, as shown in Figure 4.A1. The squared residue of such matrix A and block structure g is zero, and thus $\mathbf{x}^\top E^{(g)} \mathbf{x} = 0$ holds. However, in block structure g' satisfying $g' \neq g$, there exists at least one block of matrix A that contains two or more mutually different values, unless g' is a refinement of g , which results in $\mathbf{x}^\top E^{(g')} \mathbf{x} > 0$. In case that g' is a refinement of g , by considering an observed matrix A with block structure g' instead of g , we obtain $\mathbf{x}^\top E^{(g')} \mathbf{x} = 0$ and $\mathbf{x}^\top E^{(g)} \mathbf{x} > 0$ from the similar discussion. This contradicts the assumption that $\mathbf{x}^\top E^{(g)} \mathbf{x} = \mathbf{x}^\top E^{(g')} \mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^{np}$. \square

4.B Proof of (4.23) that $\text{rank}(E^{(\hat{g})}) = np - K_0 H_0$

Proof. For any cluster memberships \hat{g} , by simultaneously switching rows and columns with the same indices, matrix $E^{(\hat{g})}$ can be transformed into matrix $\tilde{E}^{(\hat{g})}$, which is given by

$$\tilde{E}^{(\hat{g})} = \begin{bmatrix} X^{(1)} & O & \dots & \dots & O \\ O & \ddots & \ddots & \dots & \vdots \\ \vdots & \ddots & X^{[H_0(k-1)+h]} & \ddots & \vdots \\ \vdots & \dots & \ddots & \ddots & O \\ O & \dots & \dots & O & X^{(K_0 H_0)} \end{bmatrix}, \quad (4.43)$$

where

$$\begin{aligned} X^{[H_0(k-1)+h]} &\equiv \left(X_{ij}^{[H_0(k-1)+h]} \right)_{1 \leq i \leq |I_k| |J_h|, 1 \leq j \leq |I_k| |J_h|}, \\ X_{ij}^{[H_0(k-1)+h]} &= \begin{cases} 1 - \frac{1}{|I_k| |J_h|} & \text{if } i = j, \\ -\frac{1}{|I_k| |J_h|} & \text{otherwise,} \end{cases} \\ i &= 1, \dots, |I_k| |J_h|, \quad j = 1, \dots, |I_k| |J_h|, \end{aligned} \quad (4.44)$$

for all (k, h) .

4. Statistical test on the estimated bicluster structure of a relational data matrix

Column cluster structure: $g^{(2)}$

1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
1	1	1	1	1	1	2	2	2	2	3	3
4	4	4	4	4	4	5	5	5	5	6	6
4	4	4	4	4	4	5	5	5	5	6	6
4	4	4	4	4	4	5	5	5	5	6	6
4	4	4	4	4	4	5	5	5	5	6	6
4	4	4	4	4	4	5	5	5	5	6	6
4	4	4	4	4	4	5	5	5	5	6	6
7	7	7	7	7	7	8	8	8	8	9	9
7	7	7	7	7	7	8	8	8	8	9	9
7	7	7	7	7	7	8	8	8	8	9	9
7	7	7	7	7	7	8	8	8	8	9	9

Row cluster structure: $g^{(1)}$

Figure 4.A1: A data matrix A whose squared residue σ^2 is zero with block structure g .

Let $\tilde{e}_i^{[H_0(k-1)+h]}$ be the i th column of the $[H_0(k-1)+h]$ th row block of matrix $\tilde{E}^{(\hat{g})}$. For $(k, h) \neq (k', h')$, vectors $\tilde{e}_i^{[H_0(k-1)+h]}$ and $\tilde{e}_j^{[H_0(k'-1)+h']}$ are linearly independent for an arbitrary set of (i, j) .

From here, we show that within the same $[H_0(k-1)+h]$ th block, the maximum number of linearly independent columns is $|I_k||J_h| - 1$. First, from (4.44), we have

$$\sum_{i=1}^{|I_k||J_h|} \tilde{e}_i^{[H_0(k-1)+h]} = \mathbf{0}. \quad \left(\because 1 - \frac{1}{|I_k||J_h|} + (|I_k||J_h| - 1) \left(-\frac{1}{|I_k||J_h|} \right) = 0 \right) \quad (4.45)$$

Therefore, the maximum number of linearly independent columns is smaller than $|I_k||J_h|$. Next, the columns of the indices of $i = 1, \dots, |I_k||J_h| - 1$ are linearly independent, since

$$\begin{aligned} & \sum_{i=1}^{|I_k||J_h|-1} c_i \tilde{e}_i^{[H_0(k-1)+h]} = \mathbf{0} \\ \Leftrightarrow & c_1 \left(1 - \frac{1}{|I_k||J_h|} \right) + \sum_{i \neq 1} c_i \left(-\frac{1}{|I_k||J_h|} \right) = 0, \dots, \end{aligned}$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

$$\begin{aligned}
& c_{|I_k||J_h|-1} \left(1 - \frac{1}{|I_k||J_h|} \right) + \sum_{i \neq |I_k||J_h|-1} c_i \left(-\frac{1}{|I_k||J_h|} \right) = 0 \\
\iff & c_1 + \left(-\frac{1}{|I_k||J_h|} \right) \sum_{i=1}^{|I_k||J_h|-1} c_i = 0, \dots, c_{|I_k||J_h|-1} + \left(-\frac{1}{|I_k||J_h|} \right) \sum_{i=1}^{|I_k||J_h|-1} c_i = 0 \\
\iff & c_1 = c_2 = \dots = c_{|I_k||J_h|-1}, \\
& c_i + \left(-\frac{1}{|I_k||J_h|} \right) (|I_k||J_h| - 1) c_i = 0, \text{ for all } i \\
\iff & c_1 = c_2 = \dots = c_{|I_k||J_h|-1}, \quad \frac{1}{|I_k||J_h|} c_i = 0, \text{ for all } i \\
\iff & c_1 = c_2 = \dots = c_{|I_k||J_h|-1} = 0. \tag{4.46}
\end{aligned}$$

By combining the above results, the maximum number of linearly independent columns of matrix $\tilde{E}^{(\hat{g})}$ is $\sum_{k=1}^{K_0} \sum_{h=1}^{H_0} (|I_k||J_h| - 1) = np - K_0 H_0$. Since the rank of matrix $E^{(\hat{g})}$ is equal to that of matrix $\tilde{E}^{(\hat{g})}$, we finally have $\text{rank}(E^{(\hat{g})}) = np - K_0 H_0$. \square

4.C Proof that the number of mutually different patterns of block structures with exactly $K_0 \times H_0$ blocks is lower bounded by $K_0^{n-K_0} H_0^{p-H_0}$

Proof. We give a proof of the statement in the first paragraph of Section 4.3.3. To derive a lower bound for the number of mutually different patterns of block structures, let us define a **subset** $\mathcal{G}_0^{(1)}$ of all the possible patterns of row cluster indexing as a set of all the row cluster membership vectors satisfying the following two conditions.

- n rows are clustered into **exactly** K_0 clusters.
- It can be equivalently represented in the unique form of Figure 4.C1 for some $\tilde{n} \in \{K_0, \dots, n\}$. In other words, its first $(\tilde{n} - 1)$ elements contain $1, \dots, (K_0 - 1)$ in ascending order, where \tilde{n} is the minimum row index of the K_0 th cluster.

For a fixed \tilde{n} , there are $\frac{(\tilde{n}-2)!}{(\tilde{n}-K_0)!(K_0-2)!}$ possible patterns of the first $(\tilde{n} - 1)$ elements of a cluster membership vector in $\mathcal{G}_0^{(1)}$. The last $(n - \tilde{n})$ elements are arbitrary (i.e., different indexing of these elements yields mutually **not** equivalent set of row cluster memberships), which have $K_0^{n-\tilde{n}}$ patterns. Therefore, there are $\sum_{\tilde{n}=K_0}^n \frac{(\tilde{n}-2)!}{(\tilde{n}-K_0)!(K_0-2)!} K_0^{n-\tilde{n}}$ patterns of mutually different sets of row cluster memberships. From the same discussion for column cluster memberships, we obtain a lower bound for the total number κ of the patterns of

4. Statistical test on the estimated bicluster structure of a relational data matrix

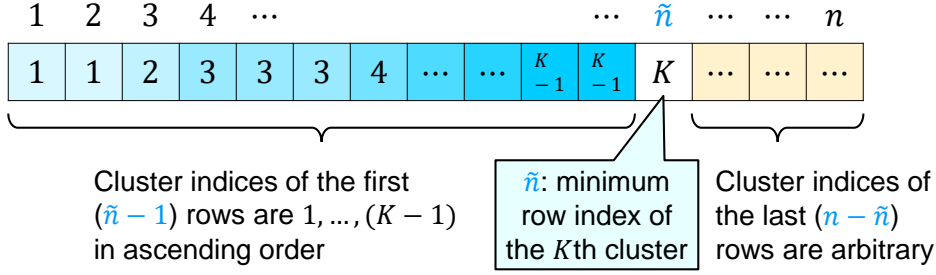


Figure 4.C1: Unique representation of row cluster indexing where n rows are clustered into **exactly** K_0 clusters. It must be noted that the set of cluster membership vectors $g^{(1)}$ that can be represented in this form is a **subset** of all the possible cluster membership vectors.

mutually different block structures:

$$\kappa \geq \left[\sum_{\tilde{n}=K_0}^n \frac{(\tilde{n} - 2)!}{(\tilde{n} - K_0)!(K_0 - 2)!} K_0^{n-\tilde{n}} \right] \left[\sum_{\tilde{p}=H_0}^p \frac{(\tilde{p} - 2)!}{(\tilde{p} - H_0)!(H_0 - 2)!} H_0^{p-\tilde{p}} \right] \geq K_0^{n-K_0} H_0^{p-H_0}, \quad (4.47)$$

which is in the exponential order of n and p for a fixed number of blocks (K_0, H_0) . \square

4.D Proof that T_E and $(\mathbf{u}_E, \mathbf{z}_E)$ are mutually independent

Proof. We give a proof of the statement in Section 4.3.1, which is used to prove (4.21). We have assumed that $\mathbf{x} \sim N(\boldsymbol{\mu}_0, \sigma_0^2 I)$ and have defined that $\mathbf{r}_E \equiv E\mathbf{x}$, $T_E = \frac{\|\mathbf{r}_E\|_2}{\sigma_0}$, $\mathbf{u}_E \equiv \frac{1}{\|\mathbf{r}_E\|_2} \mathbf{r}_E$, $\mathbf{z}_E \equiv \mathbf{x} - \mathbf{r}_E$. Note that the following equations hold:

$$\mathbf{u}_E^\top \boldsymbol{\mu}_0 = \frac{1}{\|\mathbf{r}_E\|_2} \mathbf{r}_E^\top \boldsymbol{\mu}_0 = \frac{1}{\|\mathbf{r}_E\|_2} \mathbf{x}^\top E^\top \boldsymbol{\mu}_0 = 0. \quad (4.48)$$

$$\mathbf{u}_E^\top \mathbf{u}_E = \frac{1}{\|\mathbf{r}_E\|_2^2} \mathbf{r}_E^\top \mathbf{r}_E = 1. \quad (4.49)$$

To obtain the last equation, we used the assumption that $E\boldsymbol{\mu}_0 = \mathbf{0}$.

Therefore, we have

$$\begin{aligned} p(\mathbf{x}) &= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} \|\mathbf{x} - \boldsymbol{\mu}_0\|_2^2 \right] \\ &= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} \|\mathbf{u}_E T_E \sigma_0 + \mathbf{z}_E - \boldsymbol{\mu}_0\|_2^2 \right] \end{aligned}$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

$$\begin{aligned}
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left\{ -\frac{1}{2\sigma_0^2} \left[\sigma_0^2 T_E^2 \mathbf{u}_E^\top \mathbf{u}_E + 2\sigma_0 T_E \mathbf{u}_E^\top (\mathbf{z}_E - \boldsymbol{\mu}_0) + \|\mathbf{z}_E - \boldsymbol{\mu}_0\|_2^2 \right] \right\} \\
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2} T_E^2 - \frac{1}{2\sigma_0^2} \|\mathbf{z}_E - \boldsymbol{\mu}_0\|_2^2 \right] (\because (4.20), (4.48), (4.49)) \\
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left(-\frac{1}{2} T_E^2 - \frac{1}{2\sigma_0^2} \|\mathbf{z}_E\|_2^2 \right) \exp \left[-\frac{1}{2\sigma_0^2} (-2\mathbf{z}_E^\top \boldsymbol{\mu}_0 + \|\boldsymbol{\mu}_0\|_2^2) \right]. \quad (4.50)
\end{aligned}$$

Next, we use the following Proposition 2.1 in [136]: let p_θ be a probability density function of an exponential family distribution with parameter θ , which is given by $p_\theta(\mathbf{x}) = h(\mathbf{x}) \exp\{[\boldsymbol{\eta}(\theta)]^\top T(\mathbf{x}) - \xi(\theta)\}$. Then, T is complete and sufficient for $\boldsymbol{\eta}$. From this proposition and (4.50), \mathbf{z}_E is complete and sufficient for $\boldsymbol{\mu}_0$.

We also show that (T_E, \mathbf{u}_E) are ancillary for $\boldsymbol{\mu}_0$. To prove this, we first show that T_E and \mathbf{u}_E are mutually independent. Let $\mathbf{y} \equiv \tilde{D}V\mathbf{x} \in \mathbb{R}^{np-K_0H_0}$, where V and \tilde{D} are the matrices defined in (4.15) and (4.17), respectively. From (4.18) and the fact that $\tilde{D}V\boldsymbol{\mu}_0 = \mathbf{0}$ ($\because \|\tilde{D}V\boldsymbol{\mu}_0\|_2^2 = \boldsymbol{\mu}_0^\top E\boldsymbol{\mu}_0 = 0$), we have $\mathbf{y} \sim N(\mathbf{0}, \sigma_0^2 I_{np-K_0H_0})$. Therefore, we have

$$p(\mathbf{y}) = \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np-K_0H_0}}} \exp \left(-\frac{1}{2\sigma_0^2} \|\mathbf{y}\|_2^2 \right). \quad (4.51)$$

From Proposition 2.1 in [136] and the fact that $T_E^2 = \|\mathbf{y}\|_2^2/\sigma_0^2$, T_E^2 is complete and sufficient for $\boldsymbol{\mu}_0$. Since there is a one-to-one correspondence between T_E and T_E^2 , T_E is also complete and sufficient for $\boldsymbol{\mu}_0$. Let $\tilde{\mathbf{u}}_E \equiv \mathbf{y}/\|\mathbf{y}\|_2$. Since $\tilde{\mathbf{u}}_E$ follows a uniform distribution on the surface of unit sphere and $\mathbf{u}_E = V^\top \tilde{D}^\top \tilde{\mathbf{u}}_E$, \mathbf{u}_E is ancillary for $\boldsymbol{\mu}_0$. By combining these results, T_E and \mathbf{u}_E are mutually independent from Basu's theorem [10]. Therefore, we have $p(T_E, \mathbf{u}_E) = p(T_E)p(\mathbf{u}_E)$, where $p(\cdot)$ denotes a probability density function. From the above discussion about $p(\mathbf{u}_E)$ and the fact that $T_E \sim \chi_{(np-K_0H_0)}$ (\because (4.19)), (T_E, \mathbf{u}_E) are ancillary for $\boldsymbol{\mu}_0$.

Based on the above results, \mathbf{z}_E and (T_E, \mathbf{u}_E) are independent from Basu's theorem [10]. Therefore, we have

$$p(T_E, \mathbf{u}_E, \mathbf{z}_E) = p(\mathbf{z}_E | T_E, \mathbf{u}_E) p(T_E, \mathbf{u}_E) = p(\mathbf{z}_E) p(T_E) p(\mathbf{u}_E). \quad (4.52)$$

From (4.52) and the fact that the ranges of \mathbf{z}_E , T_E , and \mathbf{u}_E do not depend on each other, we also have $p(\mathbf{u}_E, \mathbf{z}_E) = p(\mathbf{u}_E)p(\mathbf{z}_E)$ and thus

$$\begin{aligned}
p(T_E | \mathbf{u}_E, \mathbf{z}_E) &= \frac{p(T_E, \mathbf{u}_E, \mathbf{z}_E)}{p(\mathbf{u}_E, \mathbf{z}_E)} = \frac{p(\mathbf{z}_E) p(T_E) p(\mathbf{u}_E)}{p(\mathbf{u}_E, \mathbf{z}_E)} = \frac{p(\mathbf{z}_E) p(T_E) p(\mathbf{u}_E)}{p(\mathbf{u}_E) p(\mathbf{z}_E)} \\
&= p(T_E), \quad (4.53)
\end{aligned}$$

which concludes the proof. \square

4.E Sensitivity analysis with respect to the cooling schedule of simulated annealing

We conducted sensitivity analysis of the approximated version of the proposed test with respect to the cooling schedule of SA in the realizable case (i.e., $(K_0, H_0) = (K, H)$). Aside from the settings of the mean vector μ_0 and the cooling schedule of SA, we employed the same settings as in Section 4.4.3. We tried the following five cooling schedules: $T_t = 10 \times r^t$, for all $t \geq 0$, where $r = 0.99, 0.97, 0.95, 0.93, 0.91$. As for the mean vector, we used the following setting:

$$\mu_0 = 0.6 \left[\text{vec} \left(\begin{bmatrix} 0.7 & 0.55 \\ 0.5 & 0.6 \end{bmatrix} \right) - 0.5 \right] + 0.5. \quad (4.54)$$

Figures 4.E1 and 4.E2, respectively, show the histograms of the p -values of the proposed and naive approximated tests for different matrix sizes and cooling schedules r . We also plotted (i) the test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36], for the p -values of the proposed and naive tests, and (ii) the accuracy of the approximated clustering algorithm in Figures 4.E3 and 4.E4, respectively. Figure 4.E5 shows the FPR and TPR. From Figure 4.E4, we see that the accuracy of the SA algorithm got lower with the smaller value of r . As shown in Figure 4.E5, the FPR was low in all the settings, while the TPR of both the proposed and naive tests got lower with the larger value of r . A possible reason for this result is that with small r , the SA algorithm tends to output “bad” solutions (i.e., solutions that yield large squared residues) and thus both the proposed and naive tests can easily reject the null hypothesis.

4.F Application of computationally efficient biclustering algorithm for estimating the cluster memberships

The proposed approximated test based on the SA algorithm is guaranteed to converge in probability to the global minimum solution in terms of the squared residue under the conditions given in Section 4.3.3. However, this SA algorithm requires much computation time before convergence. As another option, we can use some computationally efficient biclustering algorithm for estimating the cluster memberships \hat{g} .

There have been proposed various fast biclustering algorithms [33, 91, 140]. Among these algorithms, we applied the biclustering algorithm that has been proposed by Tan and Witten [140], which is aim to minimize the loss function $\mathcal{L}(g, B; \mathbf{x})$ in (4.8). In this algorithm, to find the local optimal solution, we iteratively estimate the block-wise mean B and row and column cluster memberships, $g^{(1)}$ and $g^{(2)}$, respectively. The specific algorithm of this method is given in Algorithm 3. There is no theoretical guarantee that this algorithm converges to the global optimal solution in terms of the squared residue

4. Statistical test on the estimated bicluster structure of a relational data matrix

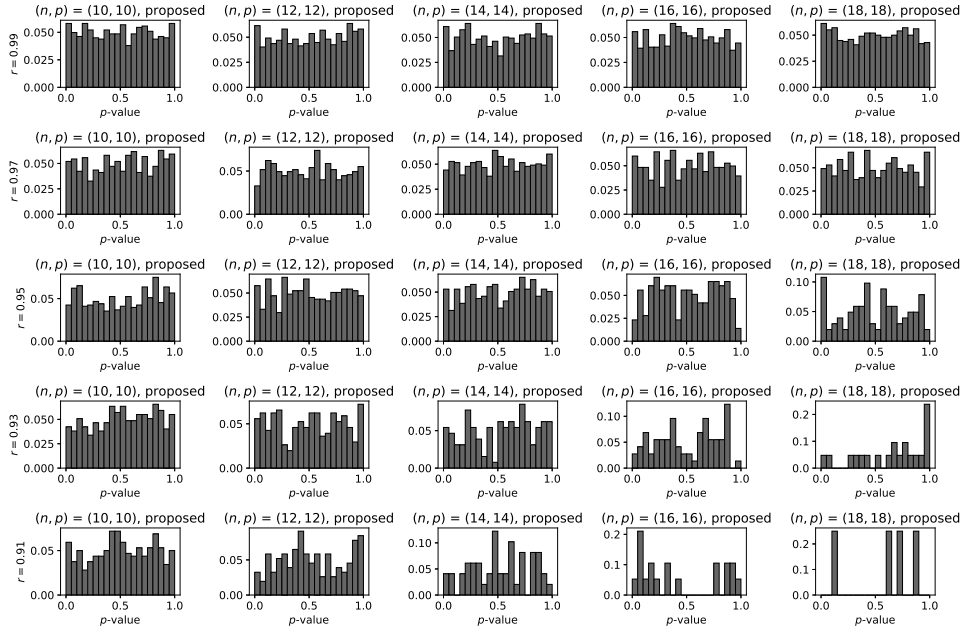


Figure 4.E1: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes and cooling schedules r , which was computed by the **approximated** version of the **proposed** test.

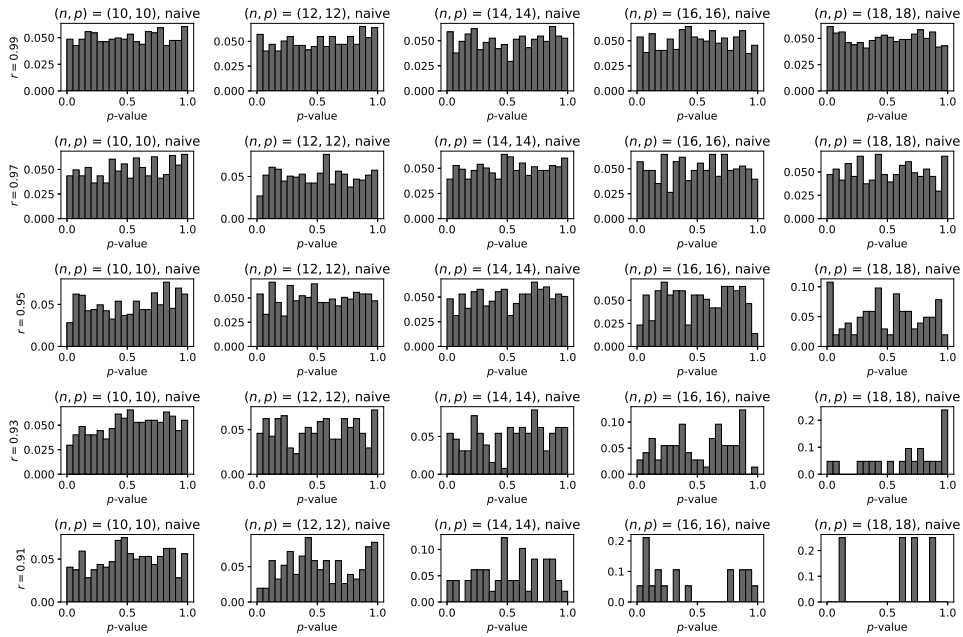


Figure 4.E2: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes and cooling schedules r , which was computed by the **approximated** version of the **naive** test (4.37).

4. Statistical test on the estimated bicluster structure of a relational data matrix

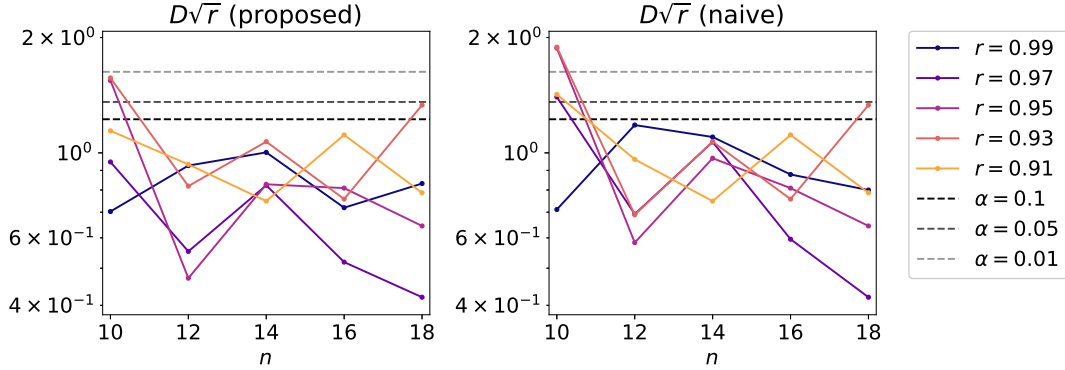


Figure 4.E3: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) **approximated** tests under the different cooling schedule settings r .

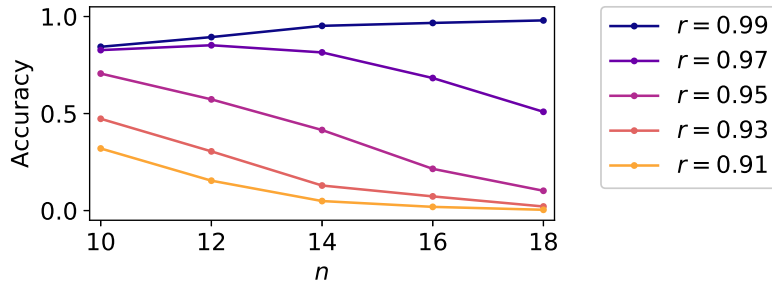


Figure 4.E4: The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and cooling schedule r , where \hat{g} is output by the **approximated** clustering algorithm in Section 4.3.3. For the experiment, we used the setting of $n = p$.

in any way, however, under the assumption that it yields a good approximation of the global optimal solution, we can use this algorithm instead of the proposed SA algorithm in Section 4.3.3 for estimating \hat{g} .

We checked the behavior of the approximated test in a realizable case when using Algorithm 3 for estimating the optimal cluster memberships \hat{g} . For finding the solution \tilde{g} of the truncation interval, we used Algorithm 2 as in the experiment in Section 4.4.3. As in Section 4.4.3, we generated data matrices and applied the approximated test. Aside from the method for estimating the cluster memberships, we used the same settings as in Section 4.4.3. This experiment was conducted on an Intel Xeon E5-2680 v3 (12 cores @ 2.50 GHz) server with 1,007 GB of RAM.

Figures 4.F1 and 4.F2, respectively, show the histograms of the p -values of the proposed and naive approximated tests. We also plotted (i) the test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36], for the p -values of the proposed and naive tests, and (ii) the accuracy of the biclustering algorithm in [140] in Figures 4.F3 and 4.F4, respectively. Figure 4.F5

4. Statistical test on the estimated bicluster structure of a relational data matrix

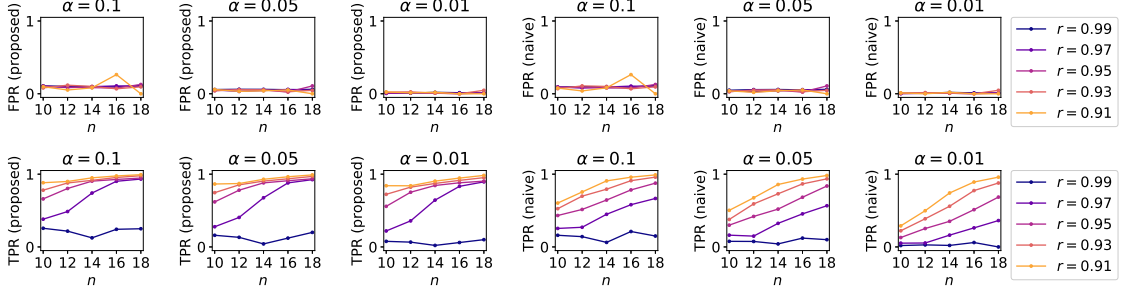


Figure 4.E5: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05,$ and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests under the different cooling schedule settings r .

shows the FPR and TPR. Finally, we plotted the computation time for each setting of mean vector μ_0 and matrix size n in Figure 4.F6. From these figures, we see that the biclustering algorithm in [140] was able to achieve accuracy comparable to or better than the proposed SA-based algorithm in less computation time.

4.G Null distribution of test statistic with unknown variance σ_0^2

We derive the null distribution of a new test statistic in case that variance σ_0^2 is unknown based on a general framework that has been proposed in [103].

Theorem 4.G.1. *Under the null hypothesis, we have*

$$T \equiv \frac{1}{c} \frac{(\|\mathbf{r}\|_2^2 - \|\mathbf{r}_1\|_2^2)}{\|\mathbf{r}_1\|_2^2} = \frac{1}{c} \frac{\|\mathbf{r}_2\|_2^2}{\|\mathbf{r}_1\|_2^2}, \quad T|\{\hat{g}, \mathbf{u}_1, \mathbf{u}_2, \mathbf{z}, \|\mathbf{r}\|_2\} \sim F_{d_1, d_2 | \hat{M}(\hat{g})}, \quad (4.55)$$

where $\|\cdot\|_2$ and $F_{d_1, d_2 | M}$, respectively, denote the Euclid norm and the truncated F distribution with parameters d_1 and d_2 and with truncation interval of M and

$$d_1 \equiv |I_1||J_1| - 1, \quad d_2 \equiv np - K_0 H_0 - |I_1||J_1| + 1, \quad c \equiv \frac{d_1}{d_2},$$

$$\mathcal{I}^{(k, h)} \equiv \{n(j-1) + i : i \in I_k, j \in J_h\},$$

$$Q^{(\hat{g})} \equiv (Q_{ij}^{(\hat{g})})_{1 \leq i \leq np, 1 \leq j \leq np}, \quad Q_{ij}^{(\hat{g})} = \begin{cases} E_{ij}^{(\hat{g})} & \text{if } [i \in \mathcal{I}^{(1,1)}] \cap [j \in \mathcal{I}^{(1,1)}], \\ 0 & \text{otherwise,} \end{cases}$$

$$\bar{Q}^{(\hat{g})} \equiv E^{(\hat{g})} - Q^{(\hat{g})},$$

$$\mathbf{r}_1 \equiv Q^{(\hat{g})} \mathbf{x}, \quad \mathbf{r}_2 \equiv \bar{Q}^{(\hat{g})} \mathbf{x}, \quad \mathbf{r} \equiv E^{(\hat{g})} \mathbf{x},$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

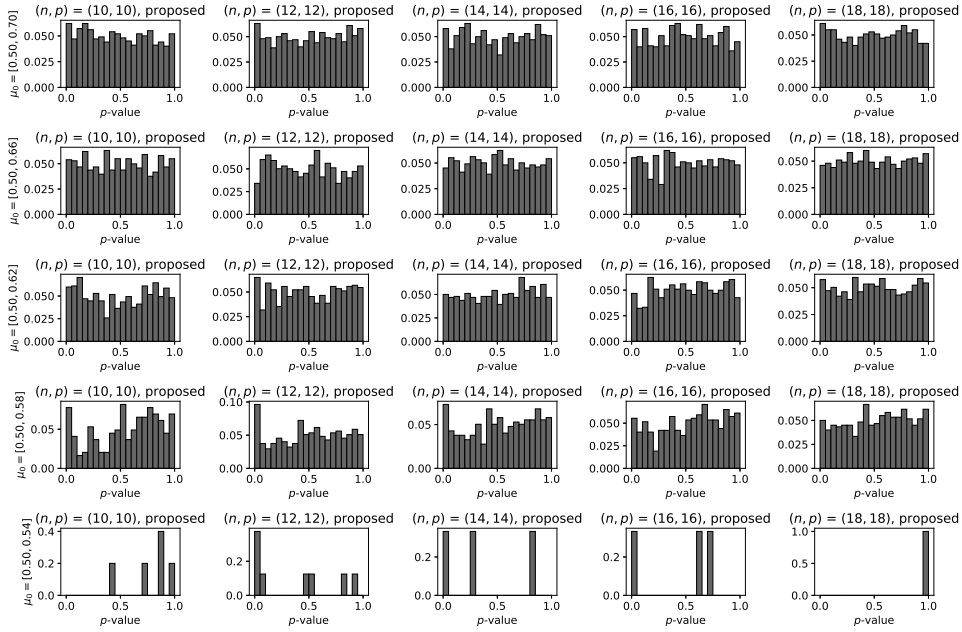


Figure 4.F1: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **approximated** version of the **proposed** test based on the biclustering algorithm in [140].

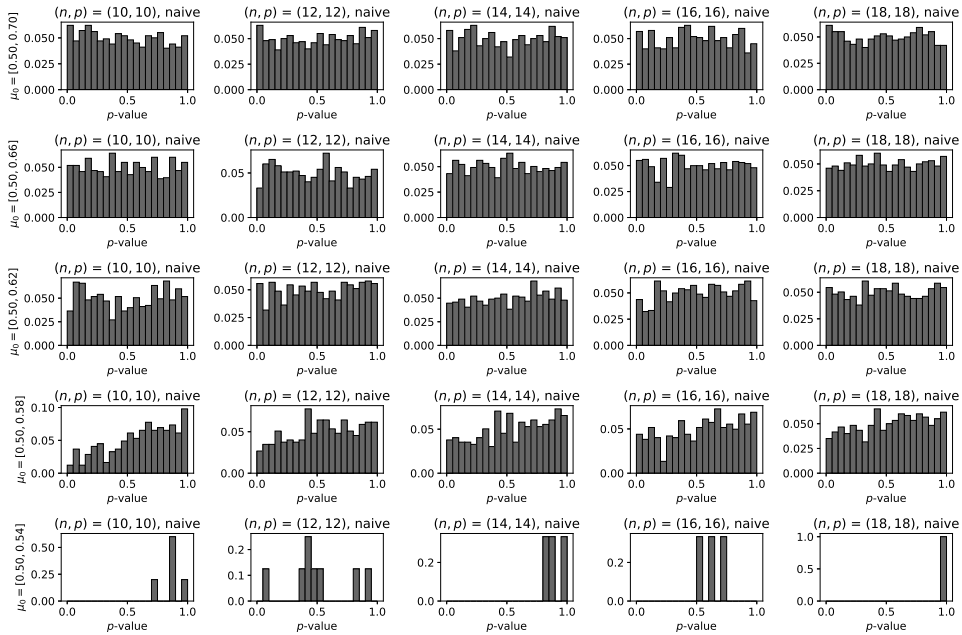


Figure 4.F2: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **approximated** version of the **naive** test (4.37) based on the biclustering algorithm in [140].

4. Statistical test on the estimated bicluster structure of a relational data matrix

Algorithm 3 Computationally efficient biclustering algorithm that has been proposed by Tan and Witten [140].

Require: A mean-centered observed matrix $\bar{A} = (\bar{A}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}$, $\bar{A}_{ij} = A_{ij} - \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p A_{ij}$.

Ensure: Approximated optimal set of cluster memberships $\hat{g} = (\hat{g}^{(1)}, \hat{g}^{(2)})$.

- 1: Define that $\hat{I}_k \equiv \{i : \hat{g}_i^{(1)} = k\}$ and $\hat{J}_h \equiv \{j : \hat{g}_j^{(2)} = h\}$.
- 2: Define initial row cluster memberships $\hat{g}^{(1)}$ by applying one-way k-means clustering to the rows of matrix \bar{A} .
- 3: Define initial column cluster memberships $\hat{g}^{(2)}$ by applying one-way k-means clustering to the columns of matrix \bar{A} .
- 4: **while true do**
- 5: $\hat{g}_0^{(1)} \leftarrow \hat{g}^{(1)}, \hat{g}_0^{(2)} \leftarrow \hat{g}^{(2)}$.
- 6: $\hat{B}_{kh} \leftarrow \frac{1}{|\hat{I}_k| |\hat{J}_h|} \sum_{i \in \hat{I}_k} \sum_{j \in \hat{J}_h} \bar{A}_{ij}$.
- 7: **for** $i = 1, \dots, n$ **do**
- 8: $\hat{g}_i^{(1)} \leftarrow \arg \min_{k \in \{1, \dots, K\}} \sum_{h=1}^H \sum_{j \in \hat{J}_h} (\bar{A}_{ij} - \hat{B}_{kh})^2$.
- 9: **end for**
- 10: $\hat{B}_{kh} \leftarrow \frac{1}{|\hat{I}_k| |\hat{J}_h|} \sum_{i \in \hat{I}_k} \sum_{j \in \hat{J}_h} \bar{A}_{ij}$.
- 11: **for** $j = 1, \dots, p$ **do**
- 12: $\hat{g}_j^{(2)} \leftarrow \arg \min_{h \in \{1, \dots, H\}} \sum_{k=1}^K \sum_{i \in \hat{I}_k} (\bar{A}_{ij} - \hat{B}_{kh})^2$.
- 13: **end for**
- 14: **if** $\hat{g}_0^{(1)} = \hat{g}^{(1)}$ and $\hat{g}_0^{(2)} = \hat{g}^{(2)}$ **then**
- 15: **break**
- 16: **end if**
- 17: **end while**

$$\begin{aligned}
 \mathbf{u}_1 &\equiv \frac{1}{\|\mathbf{r}_1\|_2} \mathbf{r}_1, & \mathbf{u}_2 &\equiv \frac{1}{\|\mathbf{r}_2\|_2} \mathbf{r}_2, \\
 \mathbf{u} &\equiv \frac{1}{\|\mathbf{r}\|_2} \mathbf{r} = \frac{1}{\sqrt{cT+1}} \mathbf{u}_1 + \sqrt{\frac{cT}{cT+1}} \mathbf{u}_2, \\
 \mathbf{z} &\equiv \mathbf{x} - \mathbf{r}, \\
 \hat{M}^{(\hat{g})} &\equiv \left\{ t \geq 0 : \hat{g} \in \hat{\mathcal{M}} \left[\|\mathbf{r}\|_2 \left(\frac{1}{\sqrt{ct+1}} \mathbf{u}_1 + \sqrt{\frac{ct}{ct+1}} \mathbf{u}_2 \right) + \mathbf{z} \right] \right\}. \tag{4.56}
 \end{aligned}$$

Proof. Let Q and \bar{Q} be fixed $np \times np$ projection matrices with the ranks of d_1 and d_2 , respectively, satisfying the following conditions:

- $Q\boldsymbol{\mu}_0 = \bar{Q}\boldsymbol{\mu}_0 = \mathbf{0}$.

4. Statistical test on the estimated bicluster structure of a relational data matrix

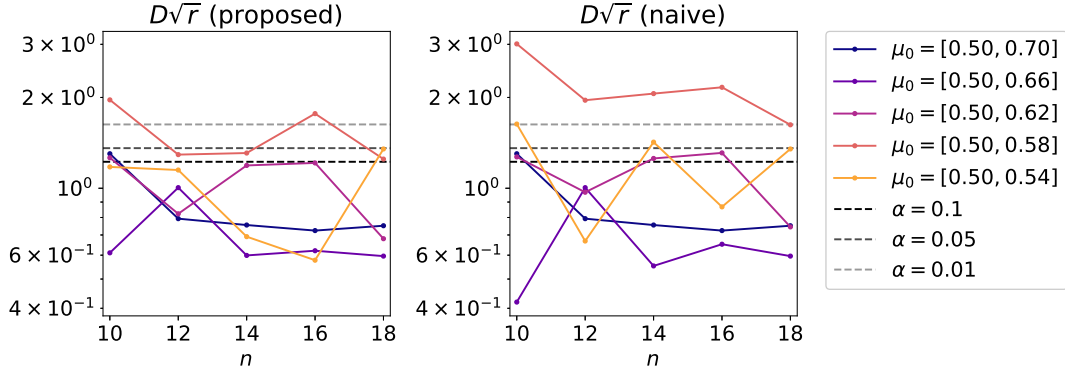


Figure 4.F3: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) **approximated** tests based on the biclustering algorithm in [140].

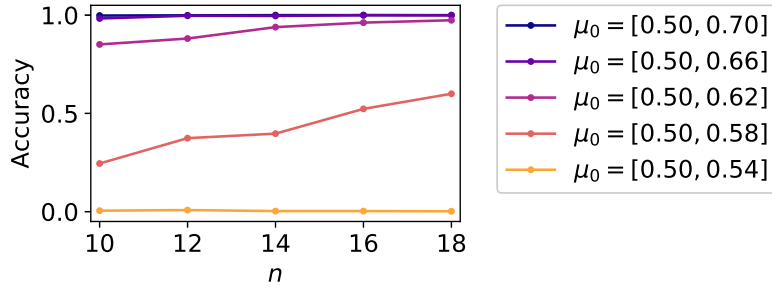


Figure 4.F4: The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the biclustering algorithm in [140]. For the experiment, we used the setting of $n = p$.

- There exists a set of row and column indices $I \subseteq \{1, \dots, np\}$ such that $Q_{ij} = 0$ if $i \notin I$ or $j \notin I$ holds and $\bar{Q}_{ij} = 0$ if $i \notin \{1, \dots, np\} \setminus I$ or $j \notin \{1, \dots, np\} \setminus I$ holds.

It must be noted that $Q\mathbf{x}$ and $\bar{Q}\mathbf{x}$ are mutually independent from the second condition.

Based on matrices Q and \bar{Q} , we use the following notations:

$$\begin{aligned}
 E &\equiv Q + \bar{Q}, \\
 \mathbf{r}_Q &\equiv Q\mathbf{x}, \quad \mathbf{r}_{\bar{Q}} \equiv \bar{Q}\mathbf{x}, \quad \mathbf{r}_E \equiv E\mathbf{x}, \\
 T_E &\equiv \frac{1}{c} \frac{(\|\mathbf{r}_E\|_2^2 - \|\mathbf{r}_Q\|_2^2)}{\|\mathbf{r}_Q\|_2^2} = \frac{1}{c} \frac{\|\mathbf{r}_{\bar{Q}}\|_2^2}{\|\mathbf{r}_Q\|_2^2}, \\
 \mathbf{u}_Q &\equiv \frac{1}{\|\mathbf{r}_Q\|_2} \mathbf{r}_Q, \quad \mathbf{u}_{\bar{Q}} \equiv \frac{1}{\|\mathbf{r}_{\bar{Q}}\|_2} \mathbf{r}_{\bar{Q}},
 \end{aligned}$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

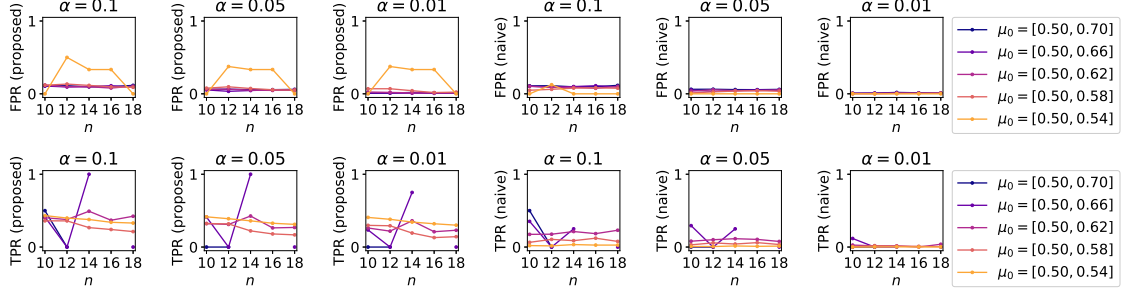


Figure 4.F5: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05,$ and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests based on the biclustering algorithm in [140].

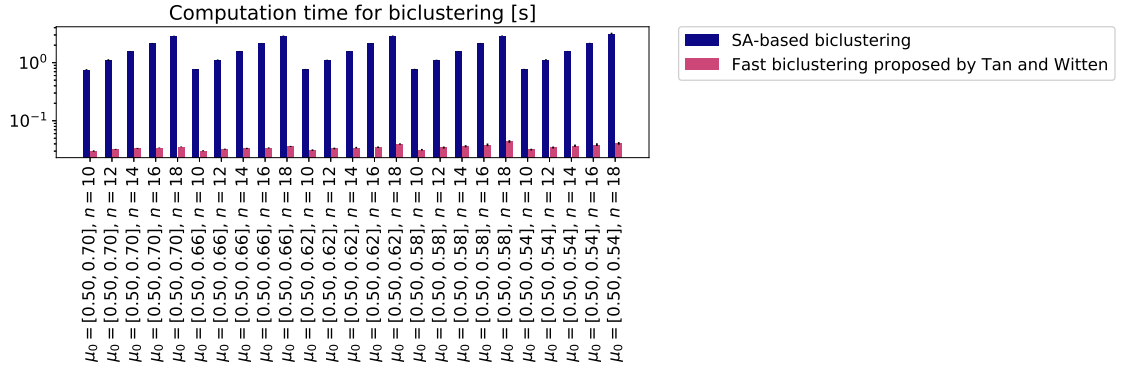


Figure 4.F6: Mean computation time for estimating the cluster memberships \hat{g} based on the proposed SA algorithm and fast biclustering algorithm in [140]. The error bars indicate the sample standard deviation of the results for 1000 trials.

$$\begin{aligned} \mathbf{u}_E &\equiv \frac{1}{\|\mathbf{r}_E\|_2} \mathbf{r}_E = \frac{1}{\sqrt{cT_E + 1}} \mathbf{u}_Q + \sqrt{\frac{cT_E}{cT_E + 1}} \mathbf{u}_{\bar{Q}}, \\ \mathbf{z}_E &\equiv \mathbf{x} - \mathbf{r}_E. \end{aligned} \quad (4.57)$$

From the similar discussion as in the proof of Theorem 3.1, we have

$$\frac{\|\mathbf{r}_Q\|_2}{\sigma_0} \sim \chi_{d_1}, \quad \frac{\|\mathbf{r}_{\bar{Q}}\|_2}{\sigma_0} \sim \chi_{d_2}. \quad (4.58)$$

Since \mathbf{r}_Q and $\mathbf{r}_{\bar{Q}}$ are mutually independent, so do $\|\mathbf{r}_Q\|_2$ and $\|\mathbf{r}_{\bar{Q}}\|_2$. By combining this fact, (4.57), and (4.58), we have

$$T_E \sim F_{d_1, d_2}. \quad (4.59)$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

Here, we show that T_E and $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2)$ are mutually independent. From the assumption, we have

$$\begin{aligned}
p(\mathbf{x}) &= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} (\mathbf{x} - \boldsymbol{\mu}_0)^\top (\mathbf{x} - \boldsymbol{\mu}_0) \right] \\
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} (\|\mathbf{r}_E\|_2 \mathbf{u}_E + \mathbf{z}_E - \boldsymbol{\mu}_0)^\top (\|\mathbf{r}_E\|_2 \mathbf{u}_E + \mathbf{z}_E - \boldsymbol{\mu}_0) \right] \\
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} (\|\mathbf{r}_E\|_2^2 \|\mathbf{u}_E\|_2^2 + \|\mathbf{z}_E - \boldsymbol{\mu}_0\|_2^2) \right] \quad (\because \mathbf{u}_E^\top \mathbf{z}_E = \mathbf{u}_E^\top \boldsymbol{\mu}_0 = 0) \\
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} (\|\mathbf{r}_E\|_2^2 + \|\mathbf{z}_E - \boldsymbol{\mu}_0\|_2^2) \right] \quad (\because \|\mathbf{u}_E\|_2 = 1) \\
&= \exp \left[-\frac{1}{2\sigma_0^2} (\|\mathbf{r}_E\|_2^2 + \|\mathbf{z}_E\|_2^2 - 2\mathbf{z}_E^\top \boldsymbol{\mu}_0) - \frac{\|\boldsymbol{\mu}_0\|_2^2}{2\sigma_0^2} - \frac{np}{2} \log(2\pi\sigma_0^2) \right]. \tag{4.60}
\end{aligned}$$

By using the notation of $\boldsymbol{\eta} \equiv \left[-\frac{1}{2\sigma_0^2} \quad \frac{1}{\sigma_0^2} \boldsymbol{\mu}_0^\top \right]^\top$, we have

$$-\frac{1}{2\sigma_0^2} (\|\mathbf{r}_E\|_2^2 + \|\mathbf{z}_E\|_2^2 - 2\mathbf{z}_E^\top \boldsymbol{\mu}_0) = \boldsymbol{\eta}^\top \left[\|\mathbf{r}_E\|_2^2 + \|\mathbf{z}_E\|_2^2 \quad \mathbf{z}_E^\top \right]^\top, \tag{4.61}$$

and thus $(\|\mathbf{r}_E\|_2^2 + \|\mathbf{z}_E\|_2^2, \mathbf{z}_E^\top)$ are complete and sufficient for $\boldsymbol{\eta}$ from Proposition 2.1 in [136]. Since there is a one-to-one correspondence between $(\mathbf{z}_E, \|\mathbf{r}_E\|_2)$ and $(\|\mathbf{r}_E\|_2^2 + \|\mathbf{z}_E\|_2^2, \mathbf{z}_E^\top)$, $(\mathbf{z}_E, \|\mathbf{r}_E\|_2)$ are also complete and sufficient for $\boldsymbol{\eta}$.

Next, we show that $(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ are ancillary for $\boldsymbol{\eta}$. To prove this, we first show that T_E and $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ are mutually independent. We use the following notations:

$$\mathbf{y}_Q \equiv \tilde{D}_Q V_Q \mathbf{x} \in \mathbb{R}^{d_1}, \quad \mathbf{y}_{\bar{Q}} \equiv \tilde{D}_{\bar{Q}} V_{\bar{Q}} \mathbf{x} \in \mathbb{R}^{d_2}, \quad \mathbf{y}_E \equiv \tilde{D}_E V_E \mathbf{x} \in \mathbb{R}^{K_0 H_0}, \tag{4.62}$$

where $Q = V_Q^\top D_Q V_Q$, $\bar{Q} = V_{\bar{Q}}^\top D_{\bar{Q}} V_{\bar{Q}}$, and $I - E = V_E^\top D_E V_E$ are singular value decompositions of matrices Q , \bar{Q} , and $I - E$, respectively, and

$$\begin{aligned}
\tilde{D}_Q &\equiv \begin{bmatrix} I_{d_1} & O_{(d_1, np-d_1)} \end{bmatrix} \in \mathbb{R}^{d_1 \times np}, \\
\tilde{D}_{\bar{Q}} &\equiv \begin{bmatrix} I_{d_2} & O_{(d_2, np-d_2)} \end{bmatrix} \in \mathbb{R}^{d_2 \times np}, \\
\tilde{D}_E &\equiv \begin{bmatrix} I_{KH} & O_{(K_0 H_0, np-K_0 H_0)} \end{bmatrix} \in \mathbb{R}^{K_0 H_0 \times np}. \tag{4.63}
\end{aligned}$$

It must be noted that we have

$$\mathbf{y}_Q \sim N(\mathbf{0}, \sigma_0^2 I_{d_1}), \quad \mathbf{y}_{\bar{Q}} \sim N(\mathbf{0}, \sigma_0^2 I_{d_2}), \quad \mathbf{y}_E \sim N(\tilde{D}_E V_E \boldsymbol{\mu}_0, \sigma_0^2 I_{K_0 H_0}). \tag{4.64}$$

From (4.60), we have

$$p(\mathbf{x}) = p(\mathbf{y}_Q, \mathbf{y}_{\bar{Q}}, \mathbf{y}_E)$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

$$\begin{aligned}
&= \frac{1}{\sqrt{(2\pi\sigma_0^2)^{np}}} \exp \left[-\frac{1}{2\sigma_0^2} (\|\mathbf{y}_Q\|_2^2 + \|\mathbf{y}_{\bar{Q}}\|_2^2 + \|\mathbf{y}_E - \tilde{D}_E V_E \boldsymbol{\mu}_0\|_2^2) \right] \\
&= p(\mathbf{y}_Q) p(\mathbf{y}_{\bar{Q}}) p(\mathbf{y}_E),
\end{aligned} \tag{4.65}$$

which results in that \mathbf{y}_Q , $\mathbf{y}_{\bar{Q}}$, and \mathbf{y}_E are independent. Since $\mathbf{r}_Q = V_Q^\top \tilde{D}_Q^\top \mathbf{y}_Q$, $\mathbf{r}_{\bar{Q}} = V_{\bar{Q}}^\top \tilde{D}_{\bar{Q}}^\top \mathbf{y}_{\bar{Q}}$, and $\mathbf{z}_E = V_E^\top \tilde{D}_E^\top \mathbf{y}_E$ hold, \mathbf{r}_Q , $\mathbf{r}_{\bar{Q}}$, and \mathbf{z}_E are also independent. Based on a similar discussion as in Appendix 4.D, $\|\mathbf{r}_Q\|_2$ and \mathbf{u}_Q are mutually independent, and so are $\|\mathbf{r}_{\bar{Q}}\|_2$ and $\mathbf{u}_{\bar{Q}}$. Therefore, we have

$$\begin{aligned}
&p(\|\mathbf{r}_Q\|_2, \mathbf{u}_Q, \|\mathbf{r}_{\bar{Q}}\|_2, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E) = p(\|\mathbf{r}_Q\|_2, \mathbf{u}_Q) p(\|\mathbf{r}_{\bar{Q}}\|_2, \mathbf{u}_{\bar{Q}}) p(\mathbf{z}_E) \\
&= p(\|\mathbf{r}_Q\|_2) p(\mathbf{u}_Q) p(\|\mathbf{r}_{\bar{Q}}\|_2) p(\mathbf{u}_{\bar{Q}}) p(\mathbf{z}_E) \\
&= p(\|\mathbf{r}_Q\|_2, \|\mathbf{r}_{\bar{Q}}\|_2) p(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}) p(\mathbf{z}_E),
\end{aligned} \tag{4.66}$$

which results in that $(\|\mathbf{r}_Q\|_2, \|\mathbf{r}_{\bar{Q}}\|_2)$ and $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ are mutually independent. Based on this fact, $T_E = \frac{1}{c} \frac{\|\mathbf{r}_{\bar{Q}}\|_2^2}{\|\mathbf{r}_Q\|_2^2}$ and $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ are mutually independent. By using this result and the fact that \mathbf{u}_Q and $\mathbf{u}_{\bar{Q}}$ are also mutually independent, we have $p(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}) = p(T_E) p(\mathbf{u}_Q) p(\mathbf{u}_{\bar{Q}})$. Based on a similar discussion as in Appendix 4.D about $p(\mathbf{u}_Q)$ and $p(\mathbf{u}_{\bar{Q}})$ and the fact that $T_E \sim F_{d_1, d_2}$ from (4.59), $(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ are ancillary for $\boldsymbol{\eta}$. Therefore, $(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ and $(\mathbf{z}_E, \|\mathbf{r}_E\|_2)$ are mutually independent from Basu's theorem [10].

By combining the above results, we have

$$\begin{aligned}
p(T_E | \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2) &= \frac{p(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2)}{p(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2)} \\
&= \frac{p(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}) p(\mathbf{z}_E, \|\mathbf{r}_E\|_2)}{p(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2)} = \frac{p(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}) p(\mathbf{z}_E, \|\mathbf{r}_E\|_2)}{p(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}) p(\mathbf{z}_E, \|\mathbf{r}_E\|_2)} \\
&= \frac{p(T_E, \mathbf{u}_Q, \mathbf{u}_{\bar{Q}})}{p(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}})} = p(T_E).
\end{aligned} \tag{4.67}$$

To derive the third equation, we used the fact that $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}})$ and $(\mathbf{z}_E, \|\mathbf{r}_E\|_2)$ are mutually independent based on a similar discussion as above. From (4.67), T_E and $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2)$ are mutually independent.

By combining the above fact and (4.59), we have

$$T_E | \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2 \sim F_{d_1, d_2}. \tag{4.68}$$

Next, we consider adding a condition of selection event of \hat{g} to the distribution of $T_E | \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2$ in (4.68). Given $(\mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2)$, the result of selection depends solely on the value of T_E , since $\mathbf{x} = \|\mathbf{r}_E\|_2 \left(\frac{1}{\sqrt{cT_E+1}} \mathbf{u}_Q + \sqrt{\frac{cT_E}{cT_E+1}} \mathbf{u}_{\bar{Q}} \right) + \mathbf{z}_E$ holds. Therefore, adding the selection condition to (4.68) corresponds to truncation of T_E to the region where $\hat{\mathcal{M}} \left[\|\mathbf{r}_E\|_2 \left(\frac{1}{\sqrt{cT_E+1}} \mathbf{u}_Q + \sqrt{\frac{cT_E}{cT_E+1}} \mathbf{u}_{\bar{Q}} \right) + \mathbf{z}_E \right] = \hat{g}$ holds:

$$T_E | \mathbf{u}_Q, \mathbf{u}_{\bar{Q}}, \mathbf{z}_E, \|\mathbf{r}_E\|_2, \hat{g} \sim F_{d_1, d_2 | \hat{\mathcal{M}}^{(\hat{g})}(E)}. \tag{4.69}$$

4. Statistical test on the estimated bicluster structure of a relational data matrix

Third, we consider replacing Q and \bar{Q} in (4.69) with $Q^{(\hat{g})}$ and $\bar{Q}^{(\hat{g})}$, which is the output by clustering algorithm \mathcal{A} based on the data vector \mathbf{x} . Based on a similar discussion as in Appendix 4.B, the matrices $Q^{(\hat{g})}$ and $\bar{Q}^{(\hat{g})}$ are also projection matrices with the ranks of d_1 and d_2 , respectively, and they satisfy the following conditions:

- $Q^{(\hat{g})}\boldsymbol{\mu}_0 = \bar{Q}^{(\hat{g})}\boldsymbol{\mu}_0 = \mathbf{0}$.
- There exists a set of row and column indices $I \subseteq \{1, \dots, np\}$ such that $Q_{ij}^{(\hat{g})} = 0$ if $i \notin I$ or $j \notin I$ holds and $\bar{Q}_{ij}^{(\hat{g})} = 0$ if $i \notin \{1, \dots, np\} \setminus I$ or $j \notin \{1, \dots, np\} \setminus I$ holds.

Since matrices $Q^{(\hat{g})}$ and $\bar{Q}^{(\hat{g})}$ depend on the data vector \mathbf{x} only through the choice of \hat{g} , under the condition that the selection result \hat{g} is given, (4.69) still holds with matrices $Q^{(\hat{g})}$ and $\bar{Q}^{(\hat{g})}$, which concludes the proof. \square

Chapter 5

Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Matrix reordering is a task to permute the rows and columns of a given observed matrix such that the resulting reordered matrix shows meaningful or interpretable structural patterns. Most existing matrix reordering techniques share the common processes of extracting some feature representations from an observed matrix in a predefined manner, and applying matrix reordering based on it. However, in some practical cases, we do not always have prior knowledge about the structural pattern of an observed matrix. To address this problem, we propose a new matrix reordering method, called deep two-way matrix reordering (DeepTMR), using a neural network model. The trained network can automatically extract nonlinear row/column features from an observed matrix, which can then be used for matrix reordering. Moreover, the proposed DeepTMR provides the denoised mean matrix of a given observed matrix as an output of the trained network. This denoised mean matrix can be used to visualize the global structure of the reordered observed matrix. We demonstrate the effectiveness of the proposed DeepTMR by applying it to both synthetic and practical datasets.

5.1 Introduction

Matrix reordering or seriation is a task to permute the rows and columns of a given observed matrix such that the resulting matrix shows meaningful or interpretable structural patterns [11, 98]. Such reordering-based matrix visualization techniques provide an overview of the various practical data matrices, including gene expression data [28, 45], document-term relationship data [15], and archaeological data [68] (e.g., the relationships between tombs and objects in Egypt [121]). In particular, we focus on the two-mode

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

two-way matrix reordering problem, where an observed matrix or *relational data matrix* $A \in \mathbb{R}^{n \times p}$ represents the relationships between two generally different objects (e.g., rows for documents and columns for words) and the permutations of rows and columns are not required to be identical, even if the row and column sizes are identical (i.e., $n = p$).

As discussed by [11], most matrix reordering techniques proposed so far share the common processes of extracting “intermediate objects” or feature representations from an observed matrix in a predefined manner, and applying matrix reordering based on the extracted intermediate objects. For instance, under *biclustering-based* methods [107, 141], which is one of the seven categories defined by [11], we assume that an observed matrix consists of homogeneous submatrices or *biclusters*, in which the entries are generated in an i.i.d. sense. Based on this assumption, we first estimate the locations (i.e., a set of row and column indices) of such biclusters and then reorder the rows and columns of the original matrix according to the estimated bicluster structure. In this example, the intermediate objects correspond to the bicluster assignments for the rows and columns.

However, in some practical cases, we do not always have prior knowledge about the structural pattern of a given observed matrix, or what input features should be used as intermediate objects. In such cases, we need to apply multiple methods and compare the results to examine which method is more suitable for analyzing the given observed matrix. Therefore, the procedure of feature extraction from an observed matrix, as well as row/column reordering, should ideally be automatically fitted to a given observed matrix.

To address this problem, we propose a new matrix reordering method, called deep two-way matrix reordering (DeepTMR), using a neural network model. The proposed DeepTMR consists of a neural network model that can be trained in an end-to-end manner and the shallower part (i.e., encoder) of the trained network can automatically extract row/column features for matrix reordering based on the given observed matrix. The expressive power of deep neural network models has been extensively studied in the literature, including the well-known universal approximation theorems [38, 52, 70]. To exploit the flexibility of neural networks for feature extraction, we transform the matrix reordering problem into a parameter estimation of the neural network, which maps row and column input features to each entry value of the observed matrix. By using an autoencoder-like neural network architecture, we train the proposed DeepTMR to extract one-dimensional row/column features from a given observed matrix, such that each entry of the observed matrix can be successfully reconstructed based on the extracted features. Then, the rows and columns are reordered based on the row/column features extracted by the trained network.

The remainder of this chapter proceeds as follows. We first review the existing matrix reordering methods and describe the differences between them and the proposed DeepTMR in Section 5.2. Then, we explain how we construct two-way matrix reordering using the proposed DeepTMR in Section 5.3. In Section 5.4, we experimentally demonstrate the effectiveness of DeepTMR by applying it to both synthetic and practical data matrices.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Finally, we discuss the results and future work directions in Section 5.5 and conclude with Section 5.6.

5.2 Related works

According to a recent survey [11], matrix reordering algorithms can be roughly classified into seven categories: *Robinsonian*, *spectral*, *dimension-reduction*, *heuristic*, *graph-theoretic*, *biclustering*, and *interactive-user-controlled*. Among them, we refer to the spectral and dimension-reduction methods, which are based on singular value decomposition (SVD) [51, 100] and multidimensional scaling (MDS) [127, 139]. These methods are particularly relevant to the proposed DeepTMR in that we assume a low-dimensional latent structure for an observed matrix. In the Robinsonian and graph-theoretic methods, the general purpose is to identify the optimal row/column orders for a given loss function [9, 24, 42, 126, 166]. However, as to obtain the global optimal solution for such a combinatorial optimization problem becomes infeasible with increasing matrix size, we need approximated algorithms for outputting local optimal solutions. For instance, finding the optimal node reordering solution for a given arbitrary graph based on bandwidth minimization or profile minimization has been shown to be NP-hard [96, 99].

Conversely, under the spectral and dimension-reduction methods (as well as the proposed DeepTMR), instead of formulating matrix reordering as a combinatorial optimization problem, we assume that an observed matrix can be well approximated by a model with a low-dimensional latent structure, estimate the parameter of this model, and interpret the estimation result as a feature of matrix reordering. By this formulation, we can avoid directly solving a combinatorial optimization problem on row/column reordering. It must be noted that biclustering-based methods are based on such a low dimensionality assumption. However, unlike the proposed DeepTMR and the dimension-reduction-based methods, they focus on detecting biclusters (i.e., a set of submatrices with coherent patterns) for an observed matrix, where the row/column orders within a bicluster are not considered in general, as also pointed out by [51].

Another advantage of the spectral and dimension-reduction methods is that some methods, including the proposed DeepTMR, can be used to extract the “denoised” mean information of a given observed matrix $A \in \mathbb{R}^{n \times p}$. Assuming that an observed matrix is generated from a statistical model, the true purpose of matrix reordering is to reveal the row/column orders of the denoised mean matrix, not to maximize the similarities between the adjacent rows and columns in the original data matrix with noise. The spectral and dimension-reduction methods address this problem, whereas the Robinsonian and graph-theoretic methods do not. In particular, in the following examples, the SVD-based method [100] (as well as the proposed DeepTMR) provides us with a denoised mean matrix of a given relational data matrix. For instance, based on the method of [100], we derive a rank-one approximation of the original matrix $A = \mathbf{r}\mathbf{c}^\top$, where $\mathbf{r} \in \mathbb{R}^n$ and $\mathbf{c} \in \mathbb{R}^p$.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

In this case, we can expect that approximated observed matrix $\mathbf{r}\mathbf{c}^\top$ preserves the global structure of the original matrix A , whereas the noise in each entry is removed. Such a denoised mean matrix can be used to visualize the global structure of an observed matrix, together with the reordered data matrix.

In the following two paragraphs, we refer to the basic matrix reordering methods based on SVD and MDS. Each of these methods is based on a specific assumption regarding the low-dimensionality of an observed matrix. The main advantages of the proposed DeepTMR for these conventional methods are as follows.

- The proposed DeepTMR can extract the low-dimensional row/column features from an observed matrix more flexibly than other methods. Unlike SVD-based methods, DeepTMR can be applied without a bilinear assumption. Moreover, unlike MDS, it does not require the specification of a distance function in advance to appropriately represent the relationships (i.e., proximity) between the pairs of rows/columns. The row/column encoder of DeepTMR, which applies a nonlinear mapping from a row/column to a one-dimensional feature, is automatically obtained by training a neural network model.
- Unlike MDS, DeepTMR can provide us with the denoised mean matrix of a given observed matrix, as well as row/column orders. Such a denoised mean matrix can be obtained as an output of the trained neural network and can be used to visualize the global structure of the reordered observed matrix.

SVD - (1) Rank-one approximation (SVD-Rank-One) Several studies have proposed utilizing SVD for matrix reordering [51, 100]. For instance, [100] have proposed to model an $n \times p$ observed matrix A with the following bilinear form:

$$\begin{aligned} \mathbf{r} &= (r_i)_{1 \leq i \leq n}, \quad \mathbf{c} = (c_j)_{1 \leq j \leq p}, \quad E = (E_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \\ A &= \mathbf{r}\mathbf{c}^\top + E, \end{aligned} \tag{5.1}$$

where \mathbf{r} and \mathbf{c} are the parameters corresponding to the rows and columns, respectively, and E is a residual matrix. Based on the above model, we estimate parameters \mathbf{r} and \mathbf{c} as follows:

$$\begin{aligned} \boldsymbol{\theta} &\equiv [r_1 \quad \cdots \quad r_n \quad c_1 \quad \cdots \quad c_p]^\top, \\ \hat{\boldsymbol{\theta}} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^{n+p}} \|A - \mathbf{r}\mathbf{c}^\top\|_F^2. \end{aligned} \tag{5.2}$$

It can be proven that the optimal solution of (5.2) is given by $\hat{\mathbf{r}} = \sqrt{\lambda_1} \mathbf{u}_1$ and $\hat{\mathbf{c}} = \sqrt{\lambda_1} \mathbf{v}_1$, where λ_1 is the largest singular value of matrix A and $\mathbf{u}_1 \in \mathbb{R}^n$ and $\mathbf{v}_1 \in \mathbb{R}^p$ are the corresponding row and column singular vectors, respectively. Therefore, the order of the estimated row and column parameters $\hat{\mathbf{r}}$ and $\hat{\mathbf{c}}$ can be respectively used for matrix reordering.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

SVD - (2) Angle between the top two singular vectors (SVD-Angle) Friendly [51] has pointed out that the structure of a given data matrix cannot always be represented sufficiently by a single principal component and has proposed to define the row/column orders using the angle between the top two singular vectors. Under this method, an observed matrix A is first mean-centered and scaled as follows:

$$\begin{aligned}\tilde{A}^{(0)} &= (\tilde{A}_{ij}^{(0)})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad \tilde{A}_{ij}^{(0)} = A_{ij} - \frac{1}{p} \sum_{j=1}^p A_{ij}, \\ \tilde{A} &= (\tilde{A}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad \tilde{A}_{ij} = \frac{\tilde{A}_{ij}^{(0)}}{\sqrt{\frac{1}{p} \sum_{j=1}^p (\tilde{A}_{ij}^{(0)})^2}}.\end{aligned}\quad (5.3)$$

Let \mathbf{u}_i be the row singular vector of scaled observed matrix \tilde{A} which corresponds to the i th largest singular value. Angle α_i between the top two row singular vectors is given by:

$$\alpha_i = \tan^{-1}(u_{i2}/u_{i1}) + \pi I[u_{i1} \leq 0], \quad (5.4)$$

where $I[\cdot]$ is an indicator function. The row order is determined by splitting angles $\{\alpha_i\}$ at the largest gap between two adjacent angles. The column order can then be defined in the same way as the row one, by replacing observed matrix \tilde{A} with transposed matrix \tilde{A}^\top .

MDS MDS is also a dimension reduction method that can be used for matrix reordering [127, 139]. Under MDS, we use a proximity matrix, each entry representing the distance between a pair of rows or columns. For instance, we can define a proximity matrix D for rows based on a given observed matrix $A \in \mathbb{R}^{n \times p}$ as follows:

$$D = (D_{ii'})_{1 \leq i, i' \leq n}, \quad D_{ii'} = \left(\sum_{j=1}^p (A_{ij} - A_{i'j})^2 \right)^{\frac{1}{2}}, \quad i, i' = 1, \dots, n. \quad (5.5)$$

The purpose of MDS is to obtain a k -dimensional representation, $\tilde{A} \in \mathbb{R}^{n \times k}$, of the original observed matrix, A , based on matrix D , where $k \leq n, p$. First, we define the following matrices:

$$\begin{aligned}\tilde{D} &= (\tilde{D}_{ii'})_{1 \leq i, i' \leq n}, \quad \tilde{D}_{ii'} = D_{ii'}^2, \quad i, i' = 1, \dots, n, \\ Q &= (Q_{ii'})_{1 \leq i, i' \leq n}, \quad Q_{ii'} = 1, \quad i, i' = 1, \dots, n, \\ B &= -\frac{1}{2} (I - n^{-1}Q) \tilde{D} (I - n^{-1}Q).\end{aligned}\quad (5.6)$$

It can be easily shown that B is a semi-positive definite matrix. Let λ_i and \mathbf{v}_i be the i th largest eigenvalue of matrix B and the corresponding eigenvector, respectively. The k -dimensional representation \tilde{A} of matrix A is given by:

$$\tilde{A} = [\mathbf{v}_1 \quad \cdots \quad \mathbf{v}_k] \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_k}). \quad (5.7)$$

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

It has been proven that solution \tilde{A} minimizes the *Strain*, which is given by $\mathcal{L}(\tilde{A}) = \|\tilde{A}\tilde{A}^\top - B\|_F^2$ [20]. By setting $k = 1$, we can obtain the one-dimensional row feature of observed matrix A , which can be used to determine the row order. The column order can be defined in the same way as the row one, by replacing observed matrix A with transposed matrix A^\top .

5.3 Main results: Deep two-way matrix reordering

Given an $n \times p$ observed matrix $A \in \mathbb{R}^{n \times p}$, our purpose is to reorder the row and column indices of matrix A based on a set of row and column permutations $\pi = (\pi^{\text{row}}, \pi^{\text{column}})$ such that the resulting matrix, $A^{(\pi)}$, exhibits some structure (e.g., block structure), as shown in Figure 1.

Figure 2 shows the entire network architecture of the proposed DeepTMR. To extract the row and column features of the given matrix A , we propose a new neural network model, DeepTMR, which has an autoencoder-like architecture. Under DeepTMR, the (i, j) th entry A_{ij} of the observed matrix A is estimated based on its row and column data vectors, $\mathbf{r}^{(i)}$ and $\mathbf{c}^{(j)}$, respectively, which are given by:

$$\begin{aligned} \mathbf{r}^{(i)} &= (r_{j'}^{(i)})_{1 \leq j' \leq p}, & r_{j'}^{(i)} &= A_{ij'}, \\ \mathbf{c}^{(j)} &= (c_{i'}^{(j)})_{1 \leq i' \leq n}, & c_{i'}^{(j)} &= A_{i'j}. \end{aligned} \quad (5.8)$$

Then, from these input data vectors, the features of the i th row and the j th column, g_i and h_j , respectively, are extracted by row and column encoder networks:

$$g_i = \text{ROWENC}(\mathbf{r}^{(i)}), \quad (5.9)$$

$$h_j = \text{COLUMNENC}(\mathbf{c}^{(j)}). \quad (5.10)$$

Here, $\text{ROWENC} : \mathbb{R}^p \mapsto \mathbb{R}$ and $\text{COLUMNENC}(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$ can be implemented as arbitrary neural network architectures, provided that they have a fixed number of units m and \tilde{m} in the input and output layers, respectively, that is, $(m, \tilde{m}) = (p, 1)$ for $\text{ROWENC}(\cdot)$ and $(m, \tilde{m}) = (n, 1)$ for $\text{COLUMNENC}(\cdot)$.

From these row and column features, the (i, j) th entry, A_{ij} , is estimated using a decoder network:

$$\hat{A}_{ij} = \text{DEC}(g_i, h_j), \quad (5.11)$$

where $\text{DEC} : \mathbb{R}^2 \mapsto \mathbb{R}$ can be implemented as an arbitrary neural network architecture with two input layer units and one output layer unit.

By using mini-batch learning, the entire network is trained such that the following mean squared error with the L_2 regularization term is minimized:

$$\mathcal{L} = \frac{1}{|\mathcal{I}_t|} \sum_{(i,j) \in \mathcal{I}_t} (A_{ij} - \hat{A}_{ij})^2 + \lambda \|\mathbf{w}\|_2^2, \quad (5.12)$$

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

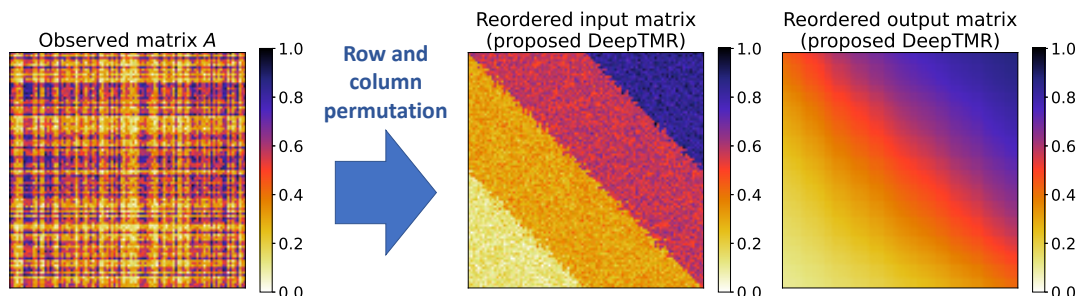


Figure 1: Matrix reordering problem. Given an observed matrix A (left), the proposed DeepTMR reorders the rows and columns of matrix A such that the reordered input matrix (center) shows a meaningful or interpretable structure. The proposed DeepTMR provides us with the denoised mean matrix of the reordered matrix (right) as the output of a trained network, as well as row/column ordering.

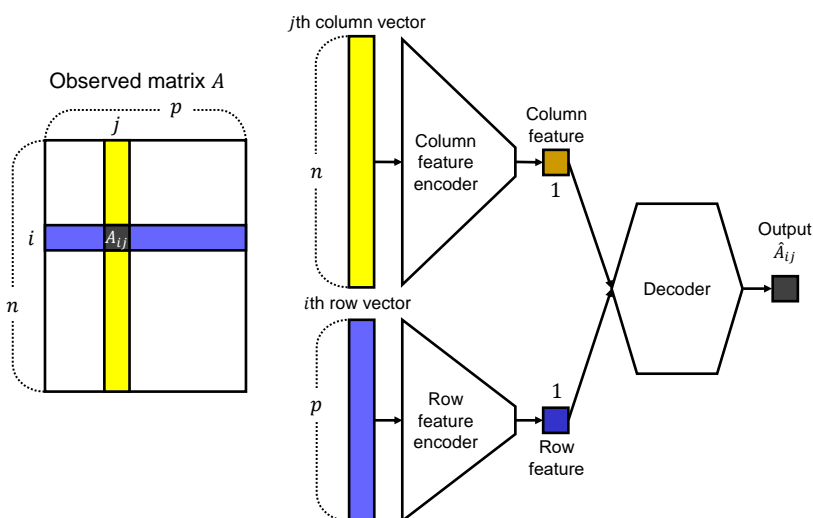


Figure 2: Model architecture of DeepTMR. Given an observed matrix A , DeepTMR is trained to reconstruct each entry A_{ij} from one-dimensional row and column features, which are extracted from the i th row and the j th column of matrix A . After training the network, we reorder the rows and columns of matrix A based on the row and column features extracted in the middle layer.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

where \mathcal{I}_t is a set of row and column indices (i, j) in a mini-batch of the t th iteration, λ is a hyperparameter, and \mathbf{w} is a vector of parameters for the entire network¹.

Finally, we define matrix $A^{(\pi)}$ with the reordered rows and columns. By using the trained row and column encoder networks, we define the following two feature vectors:

$$\begin{aligned} \mathbf{g} &= [\text{ROWENC}(\mathbf{r}^{(1)}) \quad \dots \quad \text{ROWENC}(\mathbf{r}^{(n)})]^\top, \\ \mathbf{h} &= [\text{COLUMNENC}(\mathbf{c}^{(1)}) \quad \dots \quad \text{COLUMNENC}(\mathbf{c}^{(p)})]^\top. \end{aligned} \quad (5.13)$$

Because the network has been trained to recover each entry value from only the corresponding row and column data vectors, vectors \mathbf{g} and \mathbf{h} can be expected to reflect the row and column features of the original matrix, A . Based on this discussion, we define π^{row} as a permutation of $\{1, 2, \dots, n\}$ that represents the ascending order of the entries of \mathbf{g} (i.e., $g_{\pi^{\text{row}}(1)} \leq g_{\pi^{\text{row}}(2)} \leq \dots \leq g_{\pi^{\text{row}}(n)}$ holds). Similarly, we define π^{column} as a permutation of $\{1, 2, \dots, p\}$ representing the ascending order of the entries of \mathbf{h} (i.e., $h_{\pi^{\text{column}}(1)} \leq h_{\pi^{\text{column}}(2)} \leq \dots \leq h_{\pi^{\text{column}}(p)}$ holds). Using these row and column permutations, we respectively obtain the reordered row and column features, $\mathbf{g}^{(\pi)}$ and $\mathbf{h}^{(\pi)}$, and the reordered observed and estimated matrices, $A^{(\pi)}$ and $\hat{A}^{(\pi)}$, as follows:

$$\begin{aligned} \mathbf{g}^{(\pi)} &= (g_i^{(\pi)})_{1 \leq i \leq n}, & g_i^{(\pi)} &= g_{\pi^{\text{row}}(i)}, \\ \mathbf{h}^{(\pi)} &= (h_j^{(\pi)})_{1 \leq j \leq p}, & h_j^{(\pi)} &= h_{\pi^{\text{column}}(j)}, \\ A^{(\pi)} &= (A_{ij}^{(\pi)})_{1 \leq i \leq n, 1 \leq j \leq p}, & A_{ij}^{(\pi)} &= A_{\pi^{\text{row}}(i)\pi^{\text{column}}(j)}, \\ \hat{A}^{(\pi)} &= (\hat{A}_{ij}^{(\pi)})_{1 \leq i \leq n, 1 \leq j \leq p}, & \hat{A}_{ij}^{(\pi)} &= \hat{A}_{\pi^{\text{row}}(i)\pi^{\text{column}}(j)}. \end{aligned} \quad (5.14)$$

5.4 Experiments

To verify the effectiveness of DeepTMR, we applied it to both synthetic and practical relational data matrices and plotted their latent row-column structures. For all the experiments:

- We initialized the weights and biases of the linear layers using the method described by [56]. In other words, each weight value that connects the l th and $(l+1)$ th layers is initialized based on a uniform distribution on interval $[-1/\sqrt{m^{(l)}}, 1/\sqrt{m^{(l)}}]$, where $m^{(l)}$ is the number of units in the l th layer. As for the biases, we set their initial values to zero.
- We used the Adam optimizer [78] with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 1.0 \times 10^{-8}$ for training the DeepTMR network².

¹In the experiments in Section 5.4, we define \mathbf{w} as a vector of all weights and biases in the linear layers of the encoder and decoder networks.

²As for learning rates η , we used different settings for each experiment, as shown in Table 1.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

5.4.1 Preliminary experiment using synthetic datasets

First, we generated three types of synthetic datasets with latent row and column structures, applied the DeepTMR, and checked whether we could successfully recover the latent structure of the given observed matrices. For all three models, we set the matrix size to $(n, p) = (100, 100)$.

Latent block model First, we generated a matrix based on a latent block model (LBM) [5, 58, 63]. Under an LBM, we assume that each row and column of a given matrix, $\bar{A}^{(0)} \in \mathbb{R}^{n \times p}$, belong to one of the K row and H column clusters, respectively. Let c_i and d_j be the row cluster index of the i th row and the column cluster index of the j th column of matrix $\bar{A}^{(0)}$. In this experiment, we set the number of row and column clusters at $(K, H) = (3, 3)$, and define the row and column cluster assignments as follows:

$$\begin{aligned} n^{(0)} &= \text{ceil}\left(\frac{n}{K}\right), \quad p^{(0)} = \text{ceil}\left(\frac{p}{H}\right), \\ c_1 &= \dots = c_{n^{(0)}} = 1, \quad c_{n^{(0)}+1} = \dots = c_{2n^{(0)}} = 2, \quad \dots, \quad c_{(K-1)n^{(0)}+1} = \dots = c_n = K, \\ d_1 &= \dots = d_{p^{(0)}} = 1, \quad d_{p^{(0)}+1} = \dots = d_{2p^{(0)}} = 2, \quad \dots, \quad d_{(H-1)p^{(0)}+1} = \dots = d_p = H, \end{aligned} \quad (5.15)$$

where $\text{ceil}(\cdot)$ is the ceiling function. Based on the above definitions, under an LBM, we assume that each entry of matrix $\bar{A}^{(0)}$ is independently generated from a block-wise identical distribution. Specifically, we generate each (i, j) th entry $\bar{A}_{ij}^{(0)}$ based on a Gaussian distribution with mean $B_{c_i d_j}$ and standard deviation σ given by:

$$B = \begin{bmatrix} 0.9 & 0.4 & 0.8 \\ 0.1 & 0.6 & 0.2 \\ 0.5 & 0.3 & 0.7 \end{bmatrix}, \quad \sigma = 0.05. \quad (5.16)$$

Striped pattern model To show that the DeepTMR can reveal a latent row-column structure that is not necessarily represented as a set of rectangular blocks, we also used the striped pattern model (SPM). An SPM is similar to an LBM, in that we assume that each entry of a given matrix $\bar{A}^{(0)} \in \mathbb{R}^{n \times p}$ belongs to one of the K clusters and it is independently generated from a cluster-wise identical distribution. However, unlike an LBM, we assume that the cluster assignments show a striped pattern rather than a regular grid one. Specifically, let c_{ij} be the cluster index of the (i, j) th entry $\bar{A}_{ij}^{(0)}$ of matrix $\bar{A}^{(0)}$. Under an SPM, the cluster assignment is given by:

$$\begin{aligned} n^{(0)} &= \text{ceil}\left(\frac{n+p}{K}\right), \\ c_{ij} &= \text{floor}\left(\frac{i+j-2}{n^{(0)}}\right) + 1, \quad i = 1, \dots, n, \quad j = 1, \dots, p. \end{aligned} \quad (5.17)$$

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Based on the above cluster assignments with a striped pattern, we generated each (i, j) th entry $\bar{A}_{ij}^{(0)}$ of matrix $\bar{A}^{(0)}$ based on a Gaussian distribution with mean b_{cij} and standard deviation σ given by:

$$\mathbf{b} = [0.9 \ 0.6 \ 0.3 \ 0.1], \quad \sigma = 0.05. \quad (5.18)$$

Gradation block model Under the above LBM and SPM, we assume that each entry is generated from a cluster-wise identical distribution. In the gradation block model (GBM), we consider a different case, where a matrix $\bar{A}^{(0)} \in \mathbb{R}^{n \times p}$ contains a block or submatrix with a gradation (i.e., continuous) pattern. Specifically, we assume that the (i, j) th entry $\bar{A}_{ij}^{(0)}$ of matrix $\bar{A}^{(0)}$ is generated from a Gaussian distribution with mean B_{ij} and standard deviation σ , being given by:

$$\begin{aligned} n^{(0)} &= \text{ceil} \left(\frac{n}{2} \right), \quad p^{(0)} = \text{ceil} \left(\frac{p}{2} \right), \\ B_{ij} &= \begin{cases} 0.1 & \text{if } (i > n^{(0)}) \cup (j > p^{(0)}), \\ \frac{0.8(j-1)}{p^{(0)}-1} + 0.1 & \text{if } (i \leq n^{(0)}) \cap (j \leq p^{(0)}), \end{cases} \quad i = 1, \dots, n, \quad j = 1, \dots, p, \\ \sigma &= 0.05. \end{aligned} \quad (5.19)$$

For all the above three models, once we generated matrix $\bar{A}^{(0)}$, we define matrix \bar{A} as follows:

$$\bar{A} = (\bar{A}_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}, \quad \bar{A}_{ij} = \frac{\bar{A}_{ij}^{(0)} - \min_{(i,j)=(1,1), \dots, (n,p)} \bar{A}_{ij}^{(0)}}{\max_{(i,j)=(1,1), \dots, (n,p)} \bar{A}_{ij}^{(0)} - \min_{(i,j)=(1,1), \dots, (n,p)} \bar{A}_{ij}^{(0)}}. \quad (5.20)$$

By definition, the maximum and minimum entries of matrix \bar{A} are one and zero, respectively. Then, we applied random permutation to the rows and columns of matrix \bar{A} to obtain observed matrix A . Finally, we applied the DeepTMR to observed matrix A and checked whether it could recover the latent row-column structure of matrix A . The hyperparameter settings for training the DeepTMR are listed in Table 1.

Figures 3, 4, and 5 show the results of the LBM, SPM, and GBM, respectively. For each figure, we plotted matrix \bar{A} ; observed matrix A ; reordered observed and estimated matrices, $A^{(\pi)}$ and $\hat{A}^{(\pi)}$, respectively; row and column feature vectors, \mathbf{g} and \mathbf{h} , respectively; and their reordered versions, $\mathbf{g}^{(\pi)}$ and $\mathbf{h}^{(\pi)}$. From the figures of reordered matrices $A^{(\pi)}$ and $\hat{A}^{(\pi)}$, we see that the DeepTMR can successfully extract the latent row-column structures (i.e., block structure, striped pattern, and gradation block structure) of given observed matrices. Particularly, the figures of the reordered output matrices, $\hat{A}^{(\pi)}$, show that the outputs of the DeepTMR network reflect the global structures of the given observed matrices. It must be noted that the order of the row and column indices in matrices

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Table 1: Experimental settings of learning rate η , number of epochs T (the total number of iterations is given by $\text{ceil}[Tnp/|\mathcal{I}|]$), regularization hyperparameter λ , number of sets of row and column indices in a mini-batch $|\mathcal{I}|$, and number of units in ROWENC, COLUMNENC, and DEC networks, $\mathbf{m}^{\text{ROWENC}}$, $\mathbf{m}^{\text{COLUMNENC}}$, and \mathbf{m}^{DEC} , respectively (from input to output).

	η	T	λ	$ \mathcal{I} $	$\mathbf{m}^{\text{ROWENC}}$	$\mathbf{m}^{\text{COLUMNENC}}$	\mathbf{m}^{DEC}
Sec. 5.4.1, LBM	1.0×10^{-2}	1×10^2	1.0×10^{-10}	2×10^2	$[p, 10, 1]$	$[n, 10, 1]$	$[2, 10, 1]$
Sec. 5.4.1, SPM							
Sec. 5.4.1, GBM		2×10^2					
Sec. 5.4.2, DGM							
Sec. 5.4.3							
Sec. 5.4.4	1×10^2	5×10^2					

$A^{(\pi)}$ and $\hat{A}^{(\pi)}$ that represents the latent structure is not always unique and, thus, it is not necessarily identical with that of original matrix A , as shown in these figures. For instance, the latent structures of the three models (i.e., LBM, SPM, and GBM) can also be represented by flipping or reversing the order of the row or column indices. Moreover, for an LBM, the arbitrary orders of the row or column clusters are permitted for representing the latent block structure.

From the figures of vectors $\mathbf{g}^{(\pi)}$ and $\mathbf{h}^{(\pi)}$, we see they capture the one-dimensional features of each row and column. For example, in the LBM case in Figure 3, the row (column) feature values in the same row (column) cluster are more similar than those in the mutually different row (column) clusters. In the GBM case in Figure 5, the row features are divided into two groups: one which contains the gradation pattern block and the other which does not. As for the column features, their values increase continuously within the gradation pattern block, whereas the remaining feature values are almost constant.

5.4.2 Comparison with existing matrix reordering methods

We also conducted a quantitative comparison between DeepTMR and the existing matrix reordering methods introduced in Section 5.2. For comparison, we chose the spectral/dimension-reduction methods based on SVD-Rank-One, SVD-Angle, and MDS, whose algorithms are described in Section 5.2. For the quantitative evaluation of these methods, we generated synthetic data matrices with true row/column orders, applied proposed and conventional methods, and compared their accuracies in matrix reordering. For simplicity, we considered the following data matrices, whose true row/column orders can be represented uniquely, except for row/column flipping.

Diagonal gradation model (DGM) We generated a matrix $\bar{A}^{(0)} \in \mathbb{R}^{n \times p}$ with the following diagonal gradation pattern, setting the matrix size at $(n, p) = (100, 100)$. We assumed

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

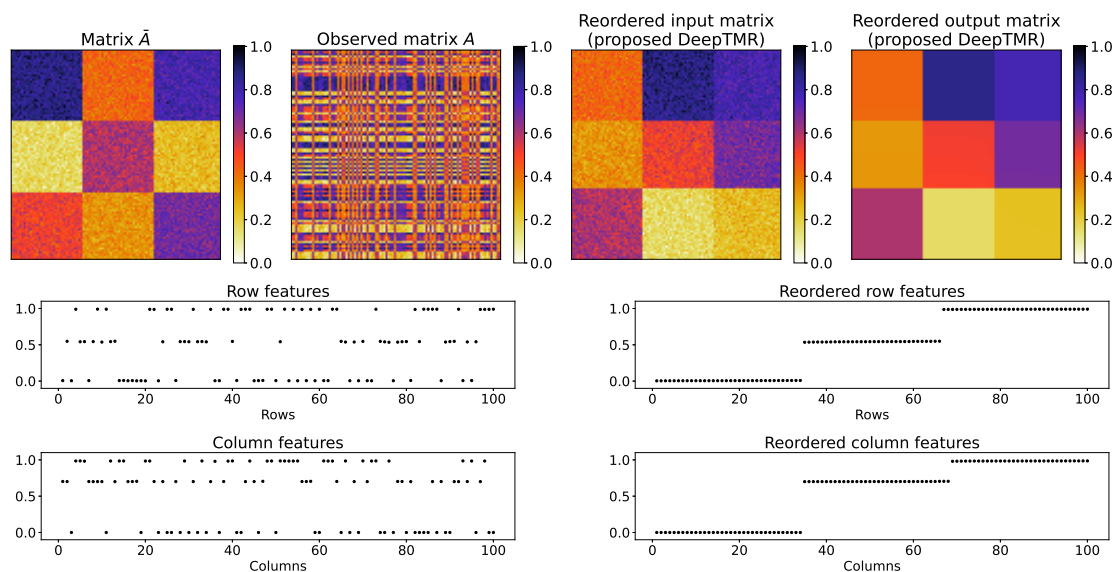


Figure 3: Results of the **LBM**. Top figures: original matrix \bar{A} , observed matrix A obtained by applying random row-column permutation to \bar{A} , reordered input matrix $A^{(\pi)}$, and reordered output matrix $\hat{A}^{(\pi)}$ (left to right). Bottom figures: Encoded row and column features \mathbf{g} and \mathbf{h} and reordered row and column features $\mathbf{g}^{(\pi)}$ and $\mathbf{h}^{(\pi)}$ (left to right).

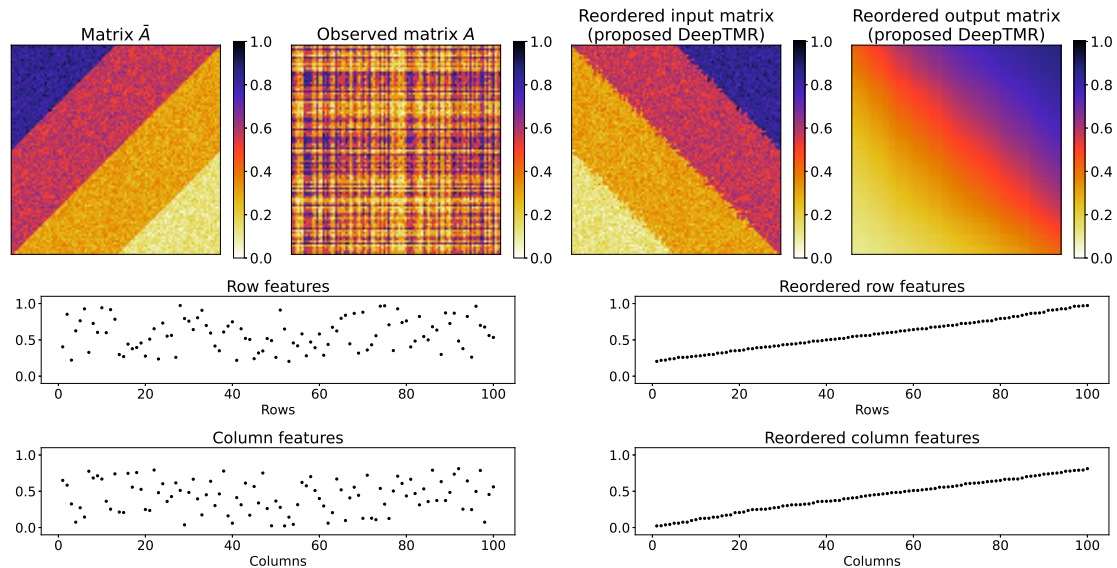


Figure 4: Results of the **SPM** for matrices \bar{A} , A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors \mathbf{g} , \mathbf{h} , $\mathbf{g}^{(\pi)}$, $\mathbf{h}^{(\pi)}$.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

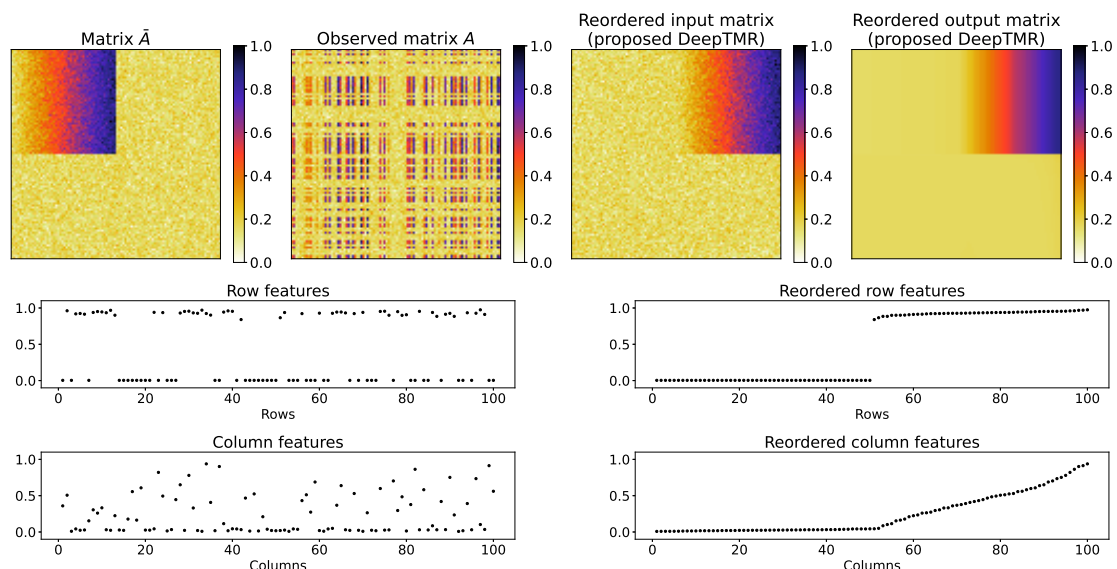


Figure 5: Results of the **GBM** for matrices \bar{A} , A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors \mathbf{g} , \mathbf{h} , $\mathbf{g}^{(\pi)}$, $\mathbf{h}^{(\pi)}$.

that the (i, j) th entry $\bar{A}_{ij}^{(0)}$ of matrix $\bar{A}^{(0)}$ is generated from a Gaussian distribution with mean B_{ij} given by:

$$B_{ij} = 0.9 - 0.8 \frac{|i - j|}{\max\{n, p\}}, \quad i = 1, \dots, n, \quad j = 1, \dots, p. \quad (5.21)$$

For the standard deviation, we tried the following 10 settings: $\sigma_t = 0.03t$ for $t = 1, \dots, 10$. As in Section 5.4.1, we defined matrix \bar{A} using matrix $\bar{A}^{(0)}$ based on (5.20), and applied a random permutation to the rows and columns of matrix \bar{A} to obtain observed matrix A . For each setting of t , we generated 10 observed matrices and applied the DeepTMR, SVD-Rank-One, SVD-Angle, and MDS. Because the training result of the DeepTMR depends on its initial parameters and the selection of the mini-batch for each iteration, for the same observed matrix, A , we trained the DeepTMR model five times and adopted the trained model with the minimum mean training loss for the last 100 iterations. The other hyperparameter settings for training the DeepTMR are listed in Table 1.

To quantitatively evaluate these methods, we computed the following matrix reordering error. Let $P \in \mathbb{R}^{n \times p}$ and $\bar{P} \in \mathbb{R}^{n \times p}$ be the population mean matrices of the reordered versions of matrix $\bar{A}^{(0)}$ (i.e., before normalization), which have the same row and column orders as matrices A and \bar{A} , respectively (i.e., $\bar{P} = B$). Let $\pi^{\text{row}(0)}$ and $\pi^{\text{column}(0)}$, respectively be the permutations of $\{1, 2, \dots, n\}$ and $\{1, 2, \dots, p\}$, which indicate the order of the rows and columns determined by each method. The flipped versions of these orders are defined as $\pi^{\text{row}(1)}$ and $\pi^{\text{column}(1)}$ (i.e., $\pi^{\text{row}(0)}(i) = \pi^{\text{row}(1)}(n - i + 1)$ and $\pi^{\text{column}(0)}(j) = \pi^{\text{column}(1)}(p - j + 1)$ for $i = 1, \dots, n$ and $j = 1, \dots, p$). Let $\bar{\pi}^{\text{row}}$

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

and $\bar{\pi}^{\text{column}}$, respectively be the order of the rows and columns that reconstruct original (correctly ordered) matrix \bar{A} . Both $\pi^{\text{row}(0)} = \bar{\pi}^{\text{row}}$ and $\pi^{\text{row}(1)} = \bar{\pi}^{\text{row}}$ indicate that the correct row ordering is obtained. Based on this fact, we redefine the row/column orders, π^{row} and π^{column} , obtained by each method as follows:

$$\begin{aligned} (\hat{k}, \hat{h}) &= \arg \min_{(k,h) \in \{(0,0), (0,1), (1,0), (1,1)\}} \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p (\bar{P}_{ij} - P_{\pi^{\text{row}(k)}(i)\pi^{\text{column}(h)}(j)})^2, \\ \pi^{\text{row}} &= \pi^{\text{row}(\hat{k})}, \quad \pi^{\text{column}} = \pi^{\text{column}(\hat{h})}. \end{aligned} \quad (5.22)$$

Finally, we define the matrix reordering error E as:

$$E = \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p (\bar{P}_{ij} - P_{\pi^{\text{row}}(i)\pi^{\text{column}}(j)})^2. \quad (5.23)$$

Figures 6 and 7 respectively show the examples of matrices \bar{A} and A with different levels of noise standard deviation σ_t , where $t = 1, \dots, 10$. Figures 8, 9, 10, and 11 respectively show the examples of the reordered observed matrix $A^{(\pi)}$ based on row/column orderings $(\pi^{\text{row}}, \pi^{\text{column}})$ obtained by DeepTMR, SVD-Rank-One, SVD-Angle, and MDS. Figure 12 shows the reordered output matrix $\hat{A}^{(\pi)}$ for the DeepTMR. From these figures, the DeepTMR and MDS can relatively successfully reorder the observed matrix, compared to the SVD-based methods. Figure 13 shows the matrix reordering error of the DeepTMR, SVD-Rank-One, SVD-Angle, and MDS. This figure shows that the DeepTMR can achieve the minimum matrix reordering error in this setting compared to the other three methods.

5.4.3 Experiment using the divorce predictors dataset

Next, we applied the DeepTMR to the divorce predictors dataset [162, 163] from the UCI Machine Learning Repository [44]. The original data matrix, $A^{(0)}$, consists of 170 rows and 54 columns, which represent the questionnaire respondents and their attributes, respectively. Each entry $A_{ij}^{(0)} \in \{0, 1, \dots, 4\}$ shows the Divorce Predictors Scale (DPS), with a higher value indicating a higher divorce risk. The meaning of the five-factor scale is as follows: 0 for “Never,” 1 for “Rarely,” 2 for “Occasionally,” 3 for “Often,” and 4 for “Always,” for Attributes 31 to 54, while they are reversed (i.e., 0 for “Always” and 4 for “Never”) for Attributes 1 to 30. The meaning of each attribute index of this dataset is provided in Appendix 5.B.

As in Section 5.4.1, we defined observed matrix A based on (5.20) by replacing $\bar{A}^{(0)}$ and \bar{A} with $A^{(0)}$ and A , respectively. Then, we applied DeepTMR to observed matrix A and checked the latent row-column structure of matrix A extracted by the DeepTMR. The hyperparameter settings for training the DeepTMR are listed in Table 1.

Figure 14 shows the results of matrices A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors \mathbf{g} , \mathbf{h} , $\mathbf{g}^{(\pi)}$, $\mathbf{h}^{(\pi)}$ for this dataset. For each row in matrices A , $A^{(\pi)}$, and $\hat{A}^{(\pi)}$, the class labels “divorced”

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

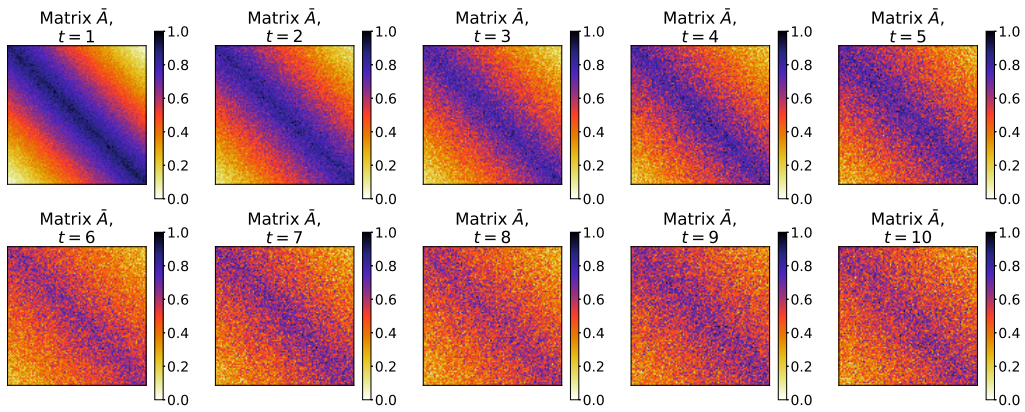


Figure 6: Examples of matrix \bar{A} for the **DGM** with different levels of noise standard deviation.

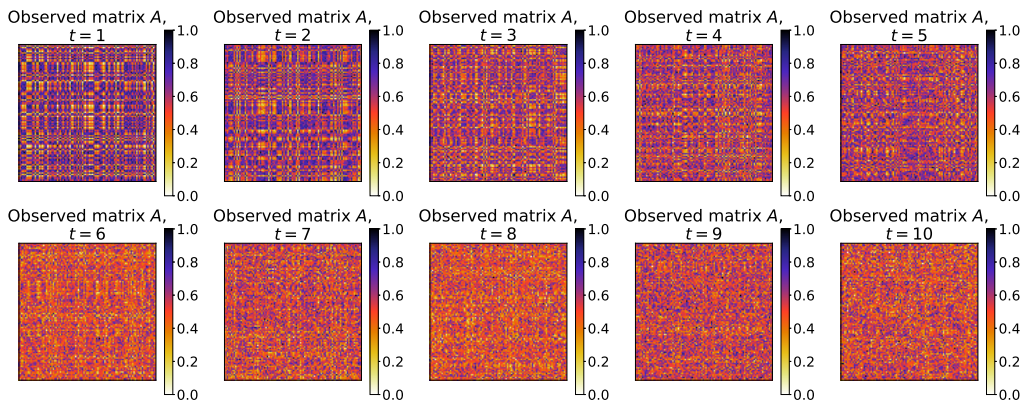


Figure 7: Examples of observed matrix A for the **DGM** with different levels of noise standard deviation.

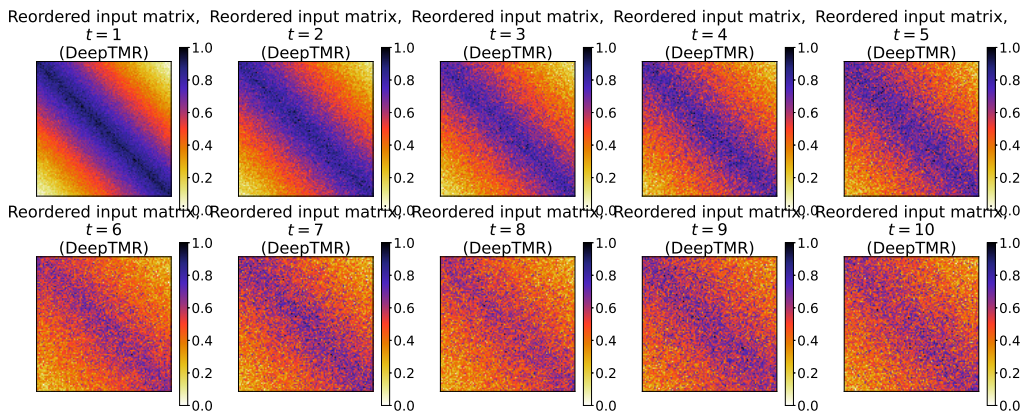


Figure 8: Examples of reordered input matrix $A^{(\pi)}$ for the **DGM** with different levels of noise standard deviation (**DeepTMR**).

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

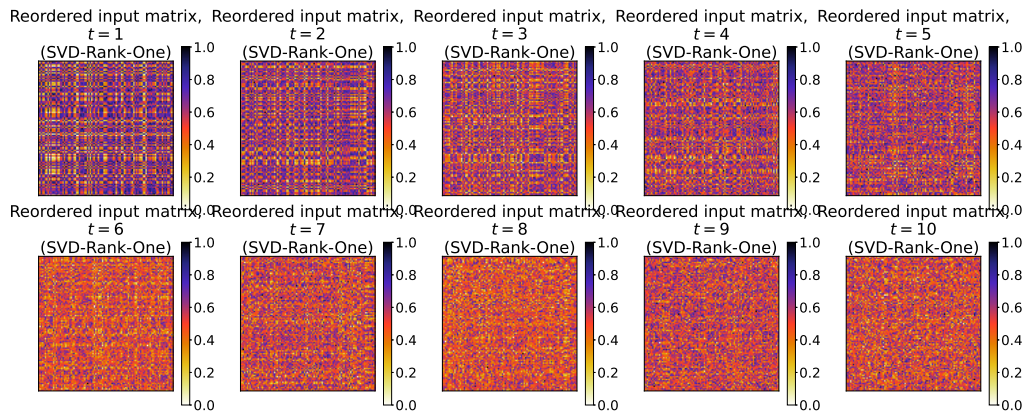


Figure 9: Examples of reordered input matrix $A^{(\pi)}$ for the **DGM** with different levels of noise standard deviation (**SVD-Rank-One**).

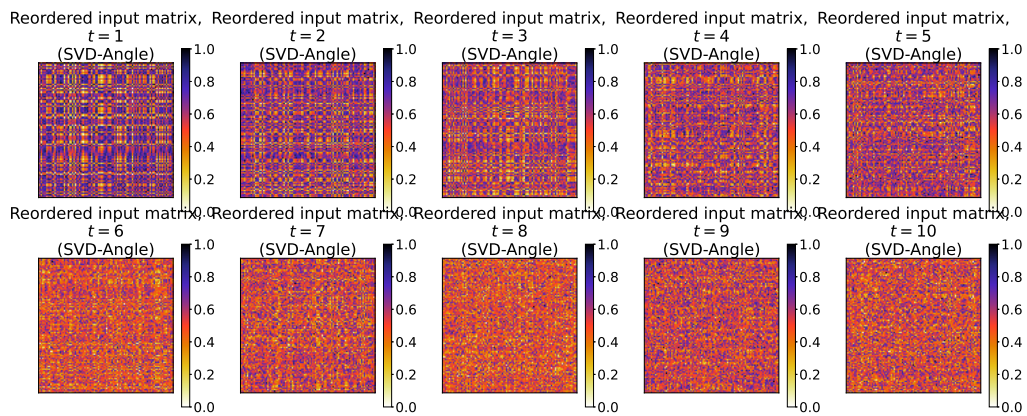


Figure 10: Examples of reordered input matrix $A^{(\pi)}$ for the **DGM** with different levels of noise standard deviation (**SVD-Angle**).

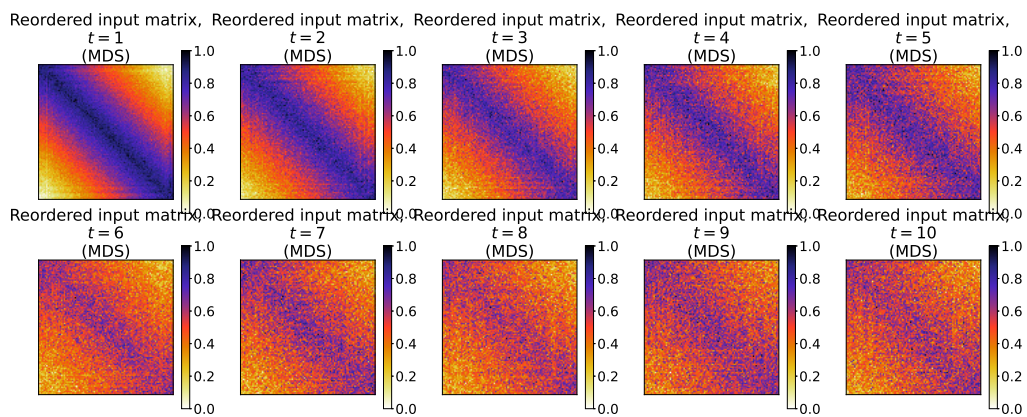


Figure 11: Examples of reordered input matrix $A^{(\pi)}$ for the **DGM** with different levels of noise standard deviation (**MDS**).

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

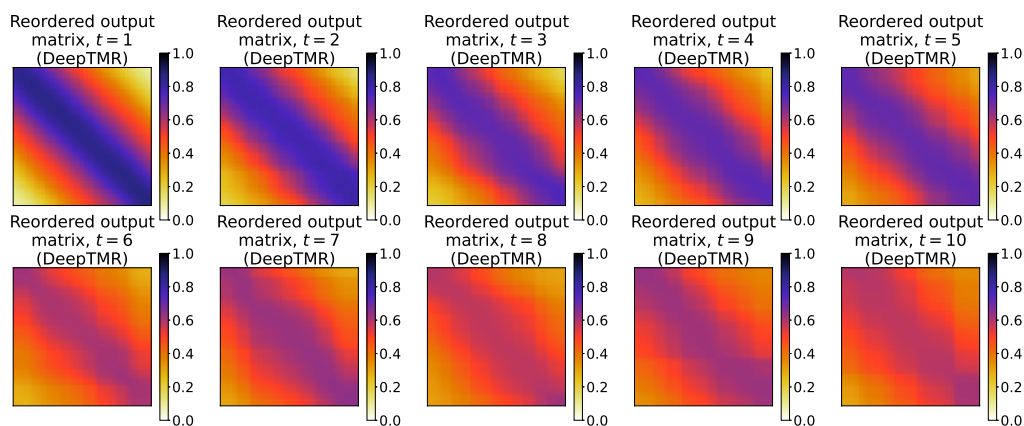


Figure 12: Examples of reordered output matrix $\hat{A}^{(\pi)}$ for the **DGM** with different levels of noise standard deviation (**DeepTMR**).

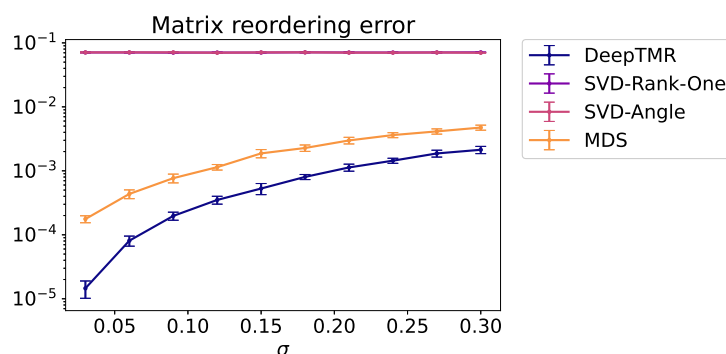


Figure 13: Matrix reordering error of the DeepTMR, SVD-Rank-One, SVD-Angle, and MDS. The error bars indicate the sample standard deviations of the results for 10 trials.

or “married” are shown in different colors on the left-hand side of the matrix. From the reordered input and output matrices $A^{(\pi)}$ and $\hat{A}^{(\pi)}$, we respectively see the latent row-column structure of the observed matrix. Roughly, the DPS takes higher values in the “divorced” rows than in the “married” ones. However, some divorced participants show relatively low DPS for some attributes (e.g., Attributes 21, 22, and 28). We also see that both the divorced and married participants show relatively high DPS for Attributes 43 and 48, whereas most participants show relatively low DPS for Attributes 6 and 7, both items referring to how to behave at home with a partner.

Figure 15 shows the results obtained with SVD-Rank-One, SVD-Angle, and MDS. From this figure, we see that the SVD-based methods (i.e., SVD-Rank-One and SVD-Angle) did not yield any meaningful structure, aside from some groups of rows with similar values. The result of MDS was similar to that of the DeepTMR, however, it did not provide

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

denoised mean information of the reordered matrix (i.e., matrix $\hat{A}^{(\pi)}$) as in Figure 14.

5.4.4 Experiment using the metropolis traffic census dataset

We also applied the DeepTMR to the metropolis traffic census dataset from e-Stat [1]. The rows and columns of the relational data matrix of this dataset represent the locations of metropolitan areas in Japan, each (i, j) th entry showing the number of people commuting (to work or school) one way from the i th location to the j th location per day. We removed the rows and columns that represent unknown locations (e.g., “unknown below Tokyo”) and the total of multiple locations (e.g., “total of three wards in central Tokyo”) from the original dataset. Let $A^{(0)} \in \mathbb{R}^{n \times p}$ be the matrix after removing the rows and columns, where $n = p = 249$. To alleviate the significant differences between entry values and consider relatively small entry values, we defined matrix $A^{(1)} \in \mathbb{R}^{n \times p}$, whose entries are given by $A_{ij}^{(1)} = \log(A_{ij}^{(0)} + 1)$ for $i = 1, \dots, n$, and $j = 1, \dots, p$.

As in Sections 5.4.1 and 5.4.3, we defined observed matrix A by replacing $\bar{A}^{(0)}$ and \bar{A} with $A^{(1)}$ and A , respectively. Then, we applied DeepTMR to observed matrix A and checked the latent row-column structure of matrix A extracted by the DeepTMR. The hyperparameter settings for training the DeepTMR are listed in Table 1.

Figures 16 and 17 respectively show the results of matrices A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors \mathbf{g} , \mathbf{h} , $\mathbf{g}^{(\pi)}$, $\mathbf{h}^{(\pi)}$ for this dataset. The correspondence of the row and column indices with locations in Figure 16 is given in Appendix 5.C. The reordered input and output matrices, $A^{(\pi)}$ and $\hat{A}^{(\pi)}$, respectively show that the number of commuting people increases from the lower right to the upper left corner of the matrices. For instance, regardless of the home location, the number of people commuting to the locations in (C21) to (C25) (e.g., Fukaya City in Saitama Prefecture, Tatebayashi City in Gunma Prefecture, and Nogi Town in Tochigi Prefecture) in matrices $A^{(\pi)}$ and $\hat{A}^{(\pi)}$ is relatively small. However, relatively many people commute to locations in (C1) (e.g., Minato-ku, Chiyoda-ku, and Shinjuku-ku in Tokyo), especially from home locations in (R1) – (R10) (e.g., Setagaya-ku, Nerima-ku, and Ota-ku in Tokyo).

Figures 18 and 19 show the results obtained with SVD-Rank-One, SVD-Angle, and MDS. As in the experiment in Section 5.4.3, the SVD-based methods (i.e., SVD-Rank-One and SVD-Angle) did not yield any meaningful structure, aside from some row/column groups with similar values. The result of MDS was similar to that of the DeepTMR except for the row and column flipping, however, it did not provide denoised mean information of the reordered matrix.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

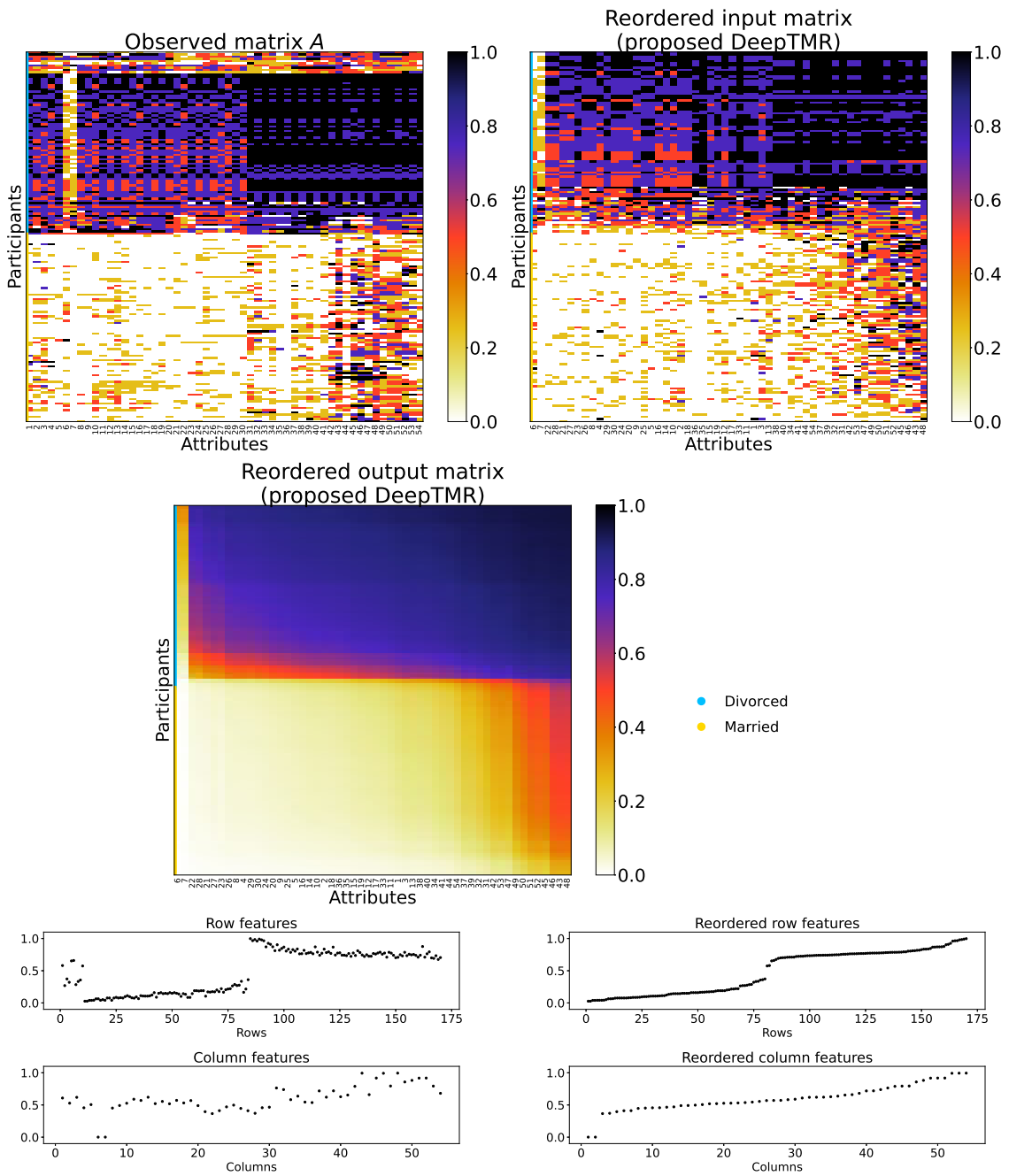


Figure 14: Results of the **divorce predictors dataset** for matrices A , $A^{(\pi)}$, $\hat{A}^{(\pi)}$, and vectors g , h , $g^{(\pi)}$, $h^{(\pi)}$.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

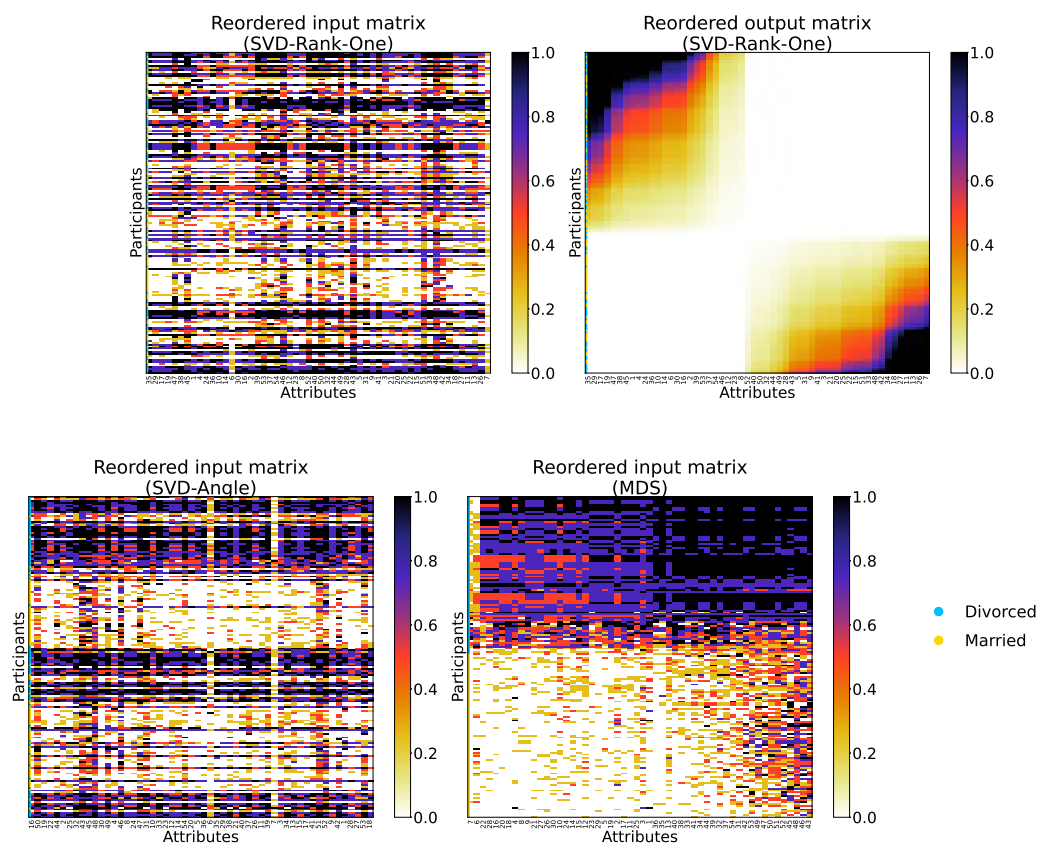


Figure 15: Results of the **divorce predictors dataset** with SVD-Rank-One, SVD-Angle, and MDS.

5.5 Discussions

Here, we discuss the results and future research directions. In the experiments in Section 5.4, we showed that the DeepTMR can successfully reorder both synthetic and practical data matrices and provide their denoised mean matrices as output. Despite its effectiveness, DeepTMR leaves room for further improvement, as described in the subsequent paragraphs.

First, one potential merit of the proposed DeepTMR compared to other spectral and dimension-reduction-based methods is that it only requires an n -dimensional column data vector and p -dimensional row data vector as input, not the entire data matrix. This suggests the possibility that, if a set of rows or columns in a data matrix increases with time, we would not have to train DeepTMR from scratch. Instead, we could only fine-tune the previously trained model with newly added data to predict orders. A main problem in realizing this is that the input dimensions of the current DeepTMR should be fixed in advance. However, to apply DeepTMR to such a time-series data matrix, we need to

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

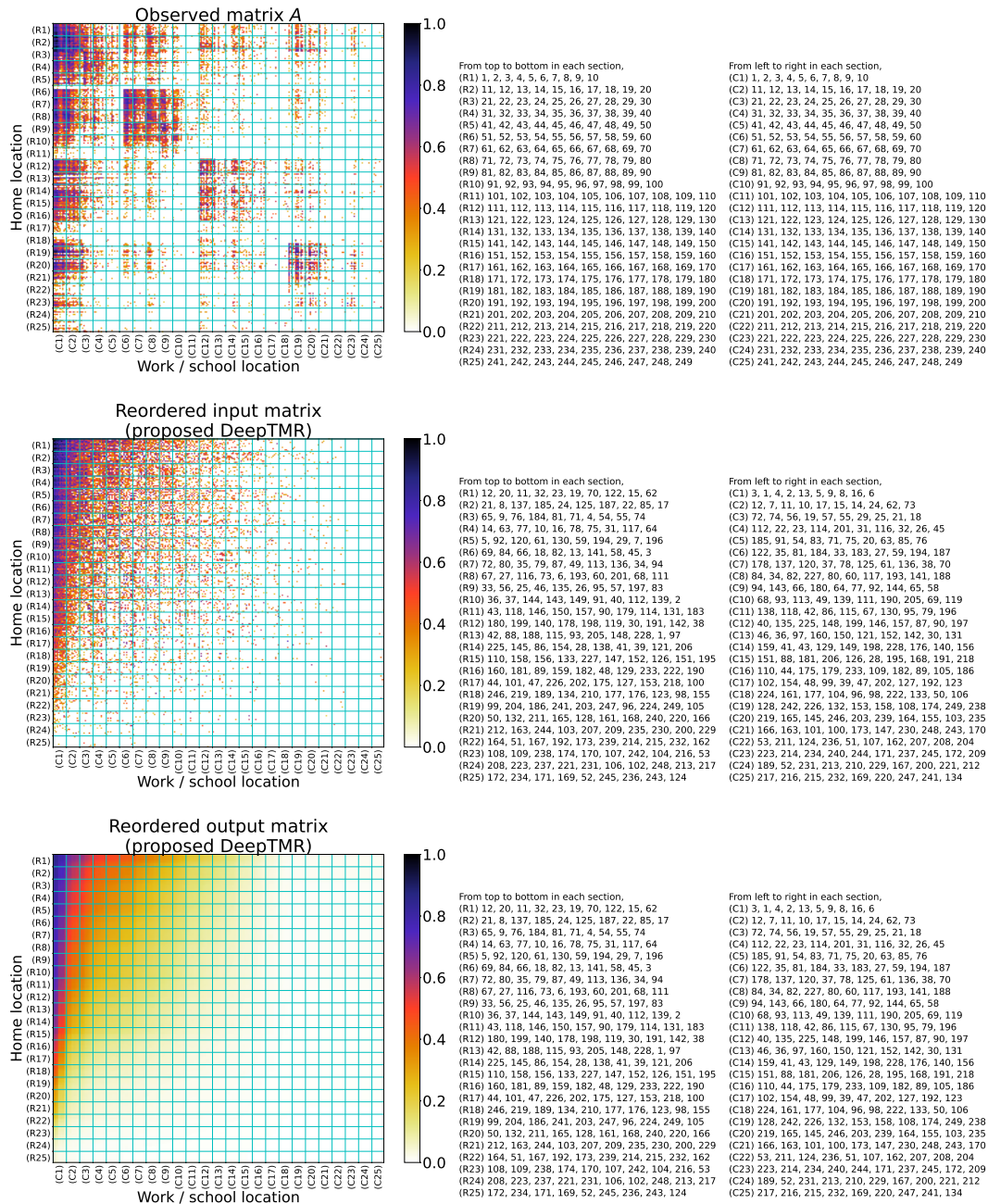


Figure 16: Results of the **metropolis traffic census dataset** for matrices A , A^π , and \hat{A}^π . For visibility, we plotted the cyan lines to show the sections between the sets of 10 rows or columns (i.e., $\{R1, \dots, R25\}$ and $\{C1, \dots, C25\}$ for rows and columns, respectively). Because the matrix size is $(n, p) = (249, 249)$, $R25$ and $C25$ contain nine rows and nine columns, respectively. The correspondence of the indices with the locations is given in Appendix 5.C.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

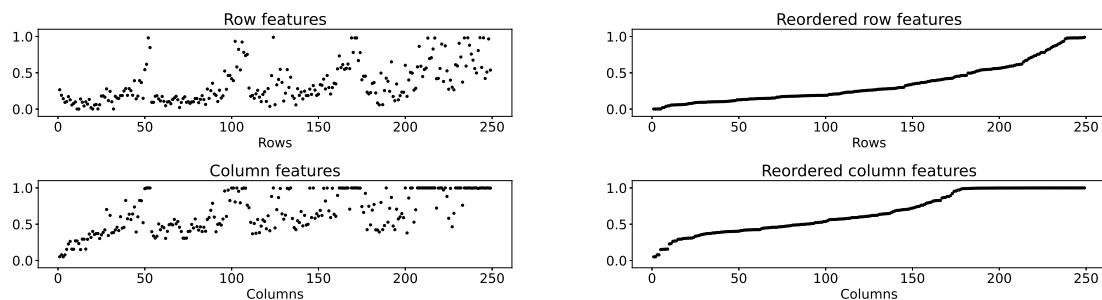


Figure 17: Results of the **metropolis traffic census dataset** for vectors g , h , $g^{(\pi)}$, and $h^{(\pi)}$.

extend the DeepTMR such that it can accept a variable-length input data vector. Such an extension is also desirable from the perspective of memory costs. For a large data matrix, a DeepTMR with an $(n + p)$ -dimensional input layer requires a large amount of memory to be stored. One possible solution to this problem is to first select randomly k rows and h columns, where $k \ll n$ and $h \ll p$, and use the selected rows and columns as inputs. In this case, we need to develop a model under a different problem setting from ours, where each entry in an input data vector does not necessarily correspond to the same row or column.

Second, another limitation of the proposed method is that the trained DeepTMR model is affected by random initialization, mini-batch selection for each iteration, and hyperparameter settings (e.g., number of units in each layer). In the experiment in Section 5.4.2, to partially alleviate this problem, we trained the neural network multiple times and chose the result with the minimum training error. However, this naïve approach increased the overall computation time. As such, it would be desirable to construct a more sophisticated model that is more robust to the effects of these settings. In particular, it is important to determine the optimal architecture of a neural network for a given data matrix. The experimental results in Section 5.4 show that the DeepTMR could successfully extract the denoised mean matrices of the input matrices with sizes ranging from 100×100 to 249×249 by using row/column encoder networks with 10 units in the middle layer. Based on these results, we expect that we do not always need to set the size of the DeepTMR network as large as the input matrix to extract the structural patterns from a data matrix.

Finally, it would be interesting to utilize additional input information for the rows and columns besides the entry values of an observed matrix. For instance, in the case of the metropolis traffic census dataset [1] in Section 5.4.4, each row or column corresponds to a specific location in Japan. If we can extend the DeepTMR to reorder a data matrix using such additional row/column information (e.g., geographical location), it would be possible to obtain a different structural pattern for the data matrix compared to those provided by the current model.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

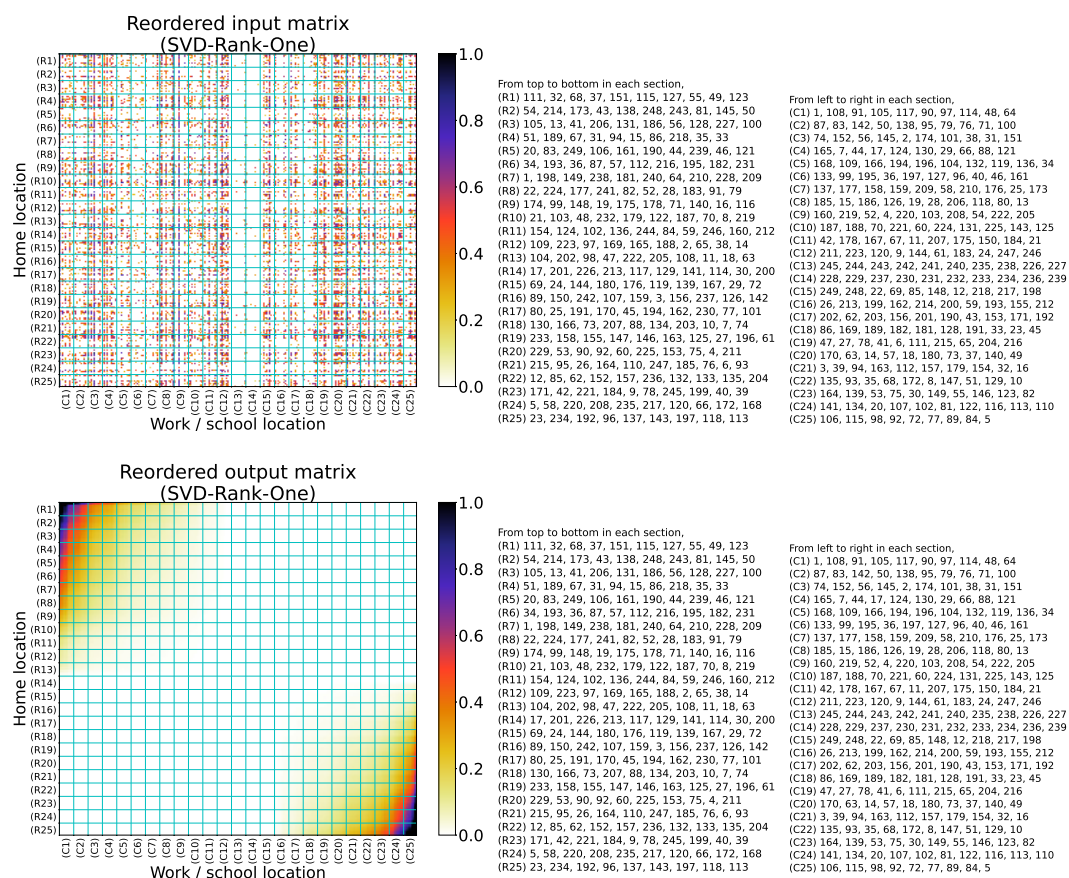


Figure 18: Results of the **metropolis** traffic census dataset with SVD-Rank-One.

5.6 Chapter conclusion

In this chapter, we proposed a new matrix reordering method, called DeepTMR, based on a neural network model. By using an autoencoder-like architecture, the proposed DeepTMR can automatically encode the row and column of an input matrix into one-dimensional nonlinear features, which can be subsequently used to determine the row and column orders. Moreover, a trained DeepTMR model provides a denoised mean matrix as output, which illustrates the global structure of the reordered input matrix. Through experiments, we showed that the proposed DeepTMR can successfully reorder the rows and columns of both synthetic and practical datasets and achieve higher accuracy in matrix reordering than the existing spectral and dimension-reduction-based matrix reordering methods based on SVD and MDS.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

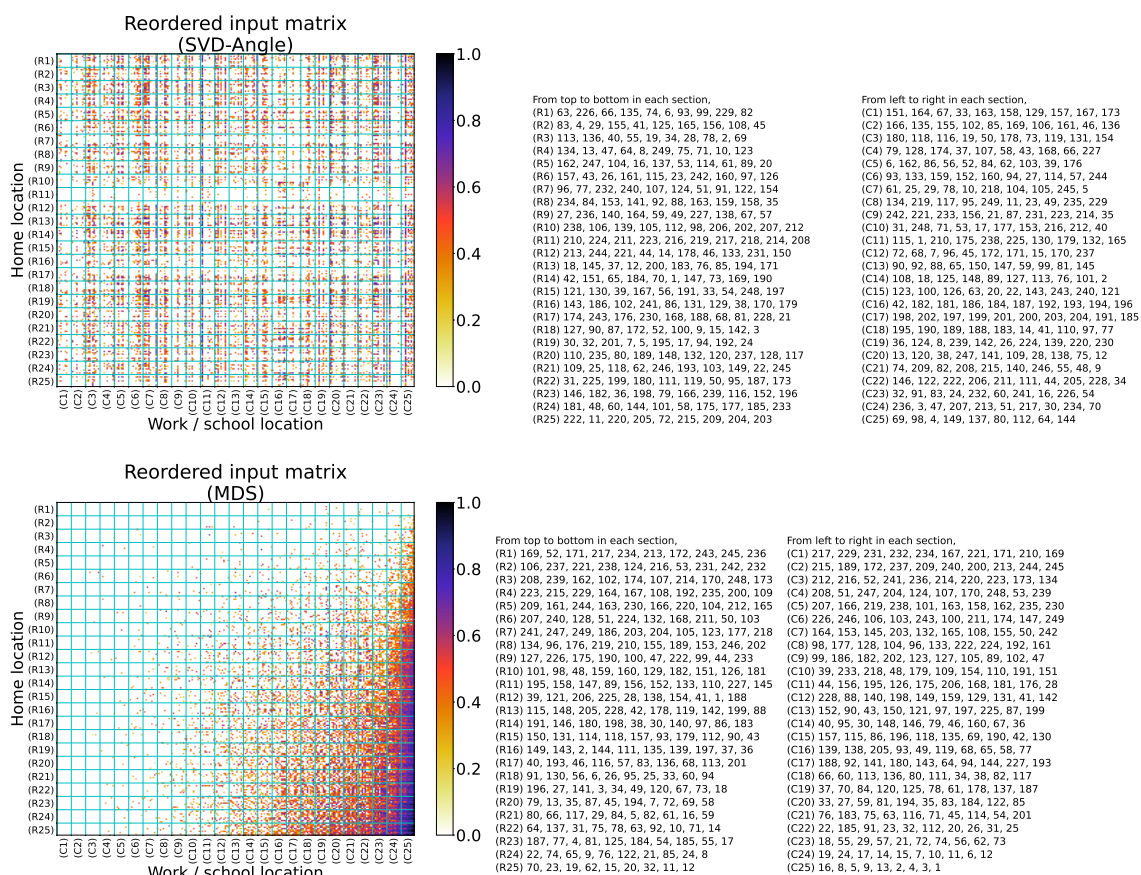


Figure 19: Results of the metropolis traffic census dataset with SVD-Angle and MDS.

5.A Application of DeepTMR to the statistical tests in Chapters 3 and 4

The proposed DeepTMR can also be used for estimating the bicluster structure of a given matrix by using its extracted row/column features. Therefore, in this section, we consider application of the DeepTMR to the statistical tests in Chapters 3 and 4. Specifically, we propose the following biclustering algorithm based on a trained DeepTMR model. First, for a given observed matrix $A^{(0)}$, we define matrix A based on (5.20) by replacing $\bar{A}^{(0)}$ and A with $A^{(0)}$ and A , respectively. Then, we train the DeepTMR model with matrix A and obtain the reordered row and column feature vectors $\mathbf{g}^{(\pi)}$ and $\mathbf{h}^{(\pi)}$ and the row and column permutations π . For a given number of row clusters K , we find the top $K - 1$ boundaries between the adjacent pairs of entries in vector $\mathbf{g}^{(\pi)}$ in terms of the absolute value of the difference. Finally, we define the row clusters by segmentation of the reordered

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Table 5.A1: Biclustering accuracy of the DeepTMR-based method and the hierarchical clustering (HC) in the settings of Gaussian, Bernoulli, and Poisson distributions.

Gaussian LBM		Bernoulli LBM		Poisson LBM	
DeepTMR	HC	DeepTMR	HC	DeepTMR	HC
0.94	1.00	0.47	0.98	0.92	1.00

row indices based on such boundaries. Based on the same procedure, we also define the column clusters by using vector $\mathbf{h}^{(\pi)}$ and a given number of column clusters H . In the experiments of the next subsections, we use this biclustering algorithm for estimating the block structure of a given matrix.

5.A.1 Application of DeepTMR to the asymptotic test on the number of biclusters in Chapter 3

First, we tried using the DeepTMR-based biclustering algorithm to compute the test statistic T of Chapter 3’s test in a realizable case. We generated 100 data matrices based on Gaussian, Bernoulli, and Poisson LBMs, estimated their bicluster structures based on DeepTMR, and computed the test statistics T . We set the matrix size to $(n, p) = (300, 225)$. Aside from the biclustering algorithm and the number of trials, we used the same settings as in Section 3.5.1. We used the same hyperparameter settings as in Section 5.4.3 for training the DeepTMR.

Table 5.A1 shows the biclustering accuracies (see Section 4.4.1 for definition) of the DeepTMR-based method and the hierarchical clustering. Figures 5.A1, 5.A2, and 5.A3, respectively, show the Q-Q plots of the test statistic T and the TW_1 distribution in the settings of Gaussian, Bernoulli, and Poisson distributions. Each plotted point corresponds to a sample of test statistic T , and the horizontal and vertical lines, respectively, show its theoretical and sample quantiles. These figures show that (1) the DeepTMR-based biclustering method could not achieve as high biclustering accuracy as hierarchical clustering and (2) the empirical distribution of the test statistic computed with hierarchical clustering was more similar to the TW_1 distribution than that computed with the DeepTMR-based biclustering method. This can be partly attributed the fact that the DeepTMR is affected by random initialization and mini-batch selection, which sometimes leads to bad local optimal solutions. To successfully apply the DeepTMR-based biclustering to the goodness-of-fit test, we need to improve its robustness to such random effect, which is beyond the scope of this dissertation.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

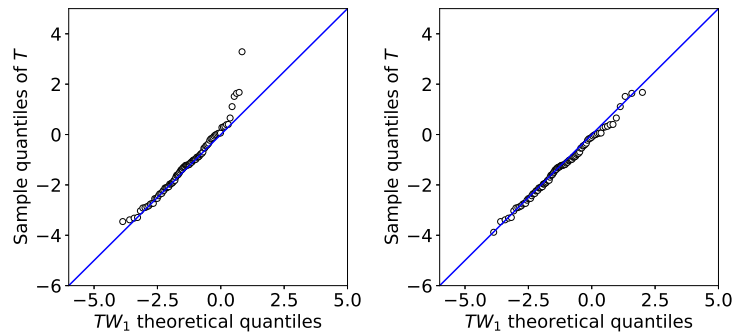


Figure 5.A1: Q-Q plot of test statistic T computed with the DeepTMR-based biclustering method (left) and hierarchical clustering (right) against the TW_1 distribution in the setting of **Gaussian case**.

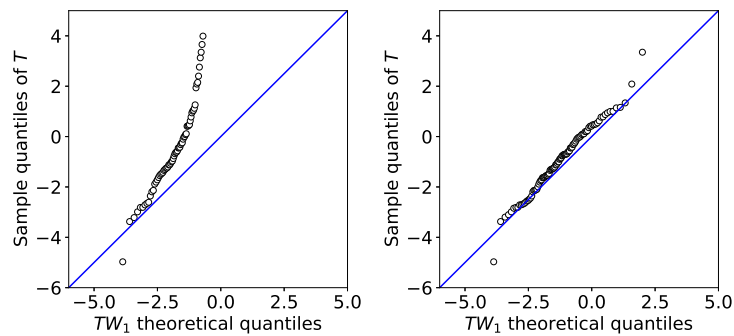


Figure 5.A2: Q-Q plot of test statistic T computed with the DeepTMR-based biclustering method (left) and hierarchical clustering (right) against the TW_1 distribution in the setting of **Bernoulli case**.

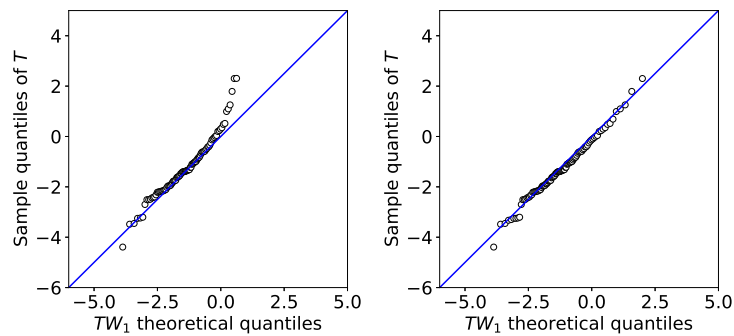


Figure 5.A3: Q-Q plot of test statistic T computed with the DeepTMR-based biclustering method (left) and hierarchical clustering (right) against the TW_1 distribution in the setting of **Poisson case**.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

5.A.2 Application of DeepTMR to the selective test on the estimated bicluster structure in Chapter 4

Although there is no theoretical guarantee that the DeepTMR-based biclustering algorithm outputs the minimum squared residue solution, as in the experiment in Section 4.F, we can use it for estimating the cluster memberships \hat{g} under the assumption that it yields a good approximation of the global minimum solution.

We generated data matrices and performed the approximated test in a realizable case when using the DeepTMR-based biclustering algorithm for estimating the optimal cluster memberships \hat{g} . For finding the solution \tilde{g} of the truncation interval, we used Algorithm 2. Aside from the biclustering algorithm, we used the same settings as in Section 4.4.3. We used the same hyperparameter settings as in Section 5.4.2 for training the DeepTMR.

Figures 5.A4 and 5.A5, respectively, show the histograms of the p -values of the proposed and naive approximated tests. We also plotted (i) the test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36], for the p -values of the proposed and naive tests, and (ii) the accuracy of the DeepTMR-based biclustering algorithm in Figures 5.A6 and 5.A7, respectively. Figures 5.A8 and 5.A9 show the FPR/TPR and the AUC score. From these figures, we see that the DeepTMR-based biclustering algorithm could not achieve as high accuracy as the SA-based algorithm. This can be partly attributed to the fact that the number of training data points (i.e., the number of entries in an observed matrix) was too small for the neural network model to be successfully trained. With regard to the AUC score, from Figure 5.A9, we see that the proposed test achieved comparable or better performance than the naive one in most settings.

5.B Correspondence of the attribute indices with meanings in the divorce predictors dataset

The meaning of each attribute index of the divorce predictors dataset [162, 163] is as follows:

1. If one of us apologizes when our discussion deteriorates, the discussion ends.
2. I know we can ignore our differences, even if things get hard sometimes.
3. When we need it, we can take our discussions with my spouse from the beginning and correct it.
4. When I discuss with my spouse, to contact him will eventually work.
5. The time I spent with my wife is special for us.
6. We don't have time at home as partners.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

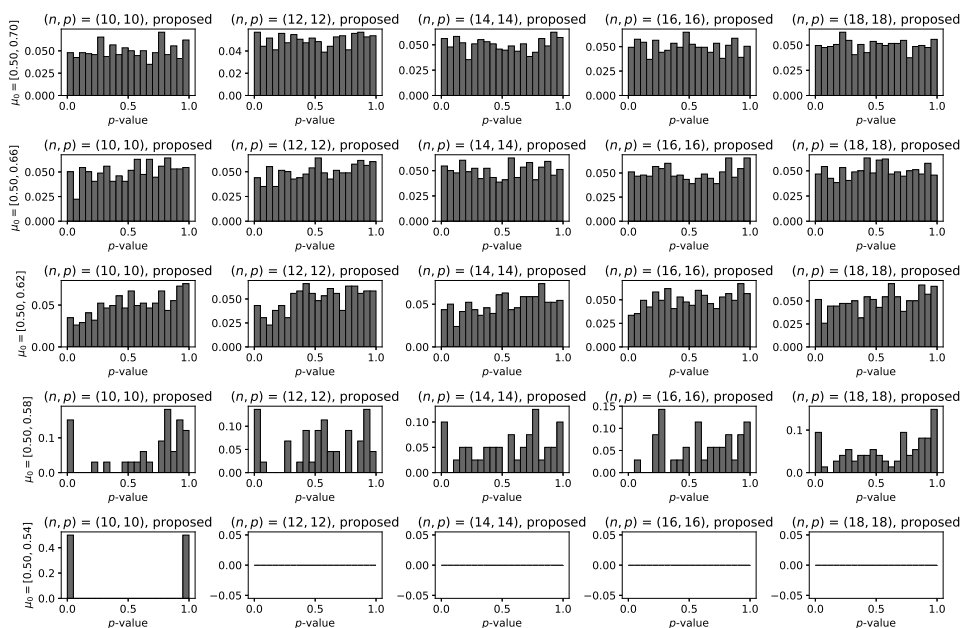


Figure 5.A4: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **approximated** version of the **proposed** test based on the biclustering algorithm using DeepTMR.

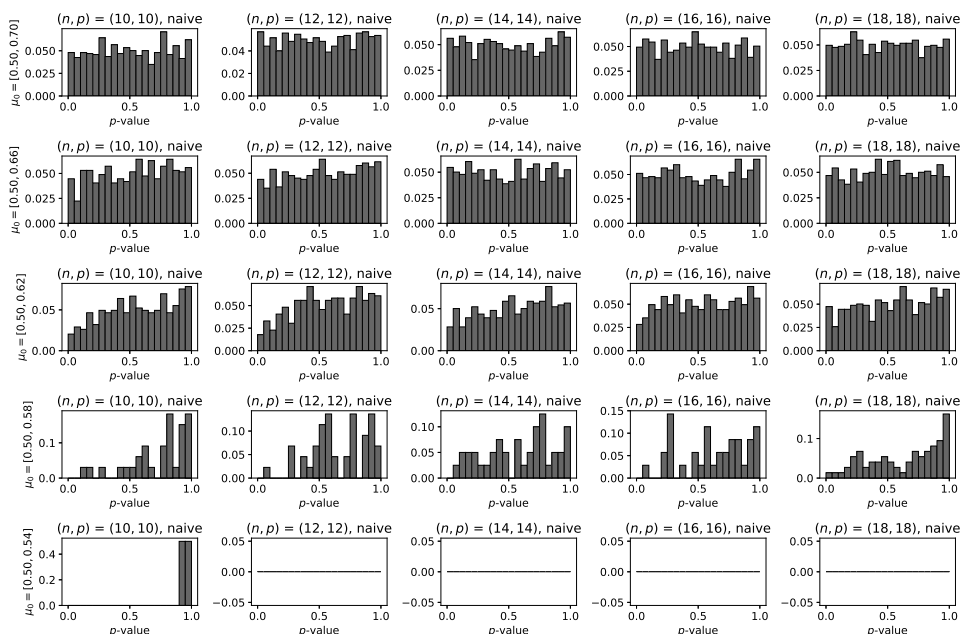


Figure 5.A5: Histograms of p -values in the null case (i.e., $\hat{g} = g^{(N)}$) for different matrix sizes, which was computed by the **approximated** version of the **naive** test (4.37) based on the biclustering algorithm using DeepTMR.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

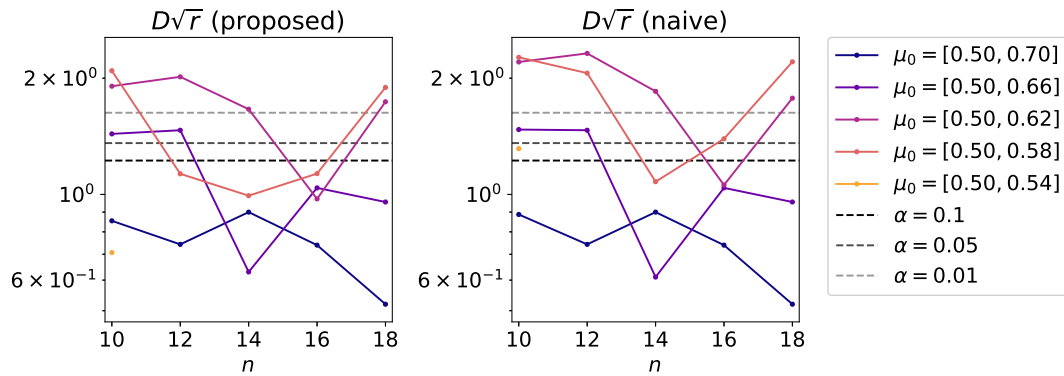


Figure 5.A6: Test statistics $D\sqrt{r}$ of the Kolmogorov-Smirnov test [36] for the p -values of the proposed (left) and naive (right) **approximated** tests based on the biclustering algorithm using DeepTMR.

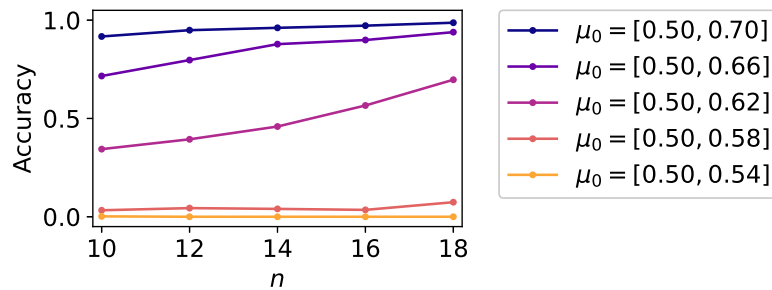


Figure 5.A7: The ratio of the number of the null cases (i.e., $\hat{g} = g^{(N)}$) for each setting of matrix size (n, p) and mean vector μ_0 , where \hat{g} is output by the biclustering algorithm using DeepTMR. For the experiment, we used the setting of $n = p$.

7. We are like two strangers who share the same environment at home rather than family.
8. I enjoy our holidays with my wife.
9. I enjoy traveling with my wife.
10. Most of our goals are common to my spouse.
11. I think that one day in the future, when I look back, I see that my spouse and I have been in harmony with each other.
12. My spouse and I have similar values in terms of personal freedom.
13. My spouse and I have similar sense of entertainment.
14. Most of our goals for people (children, friends, etc.) are the same.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

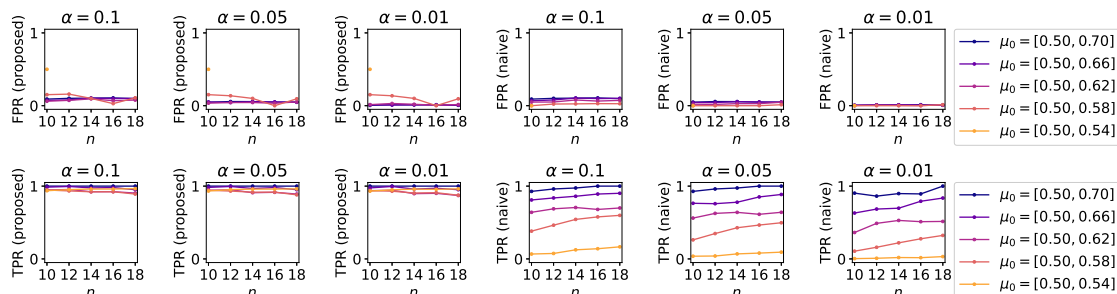


Figure 5.A8: FPR and TPR in the realizable case with different significance rates (e.g., $\alpha = 0.1, 0.05$, and 0.01), for the **approximated** version of the proposed (left) and naive (right) statistical tests based on the biclustering algorithm using DeepTMR.

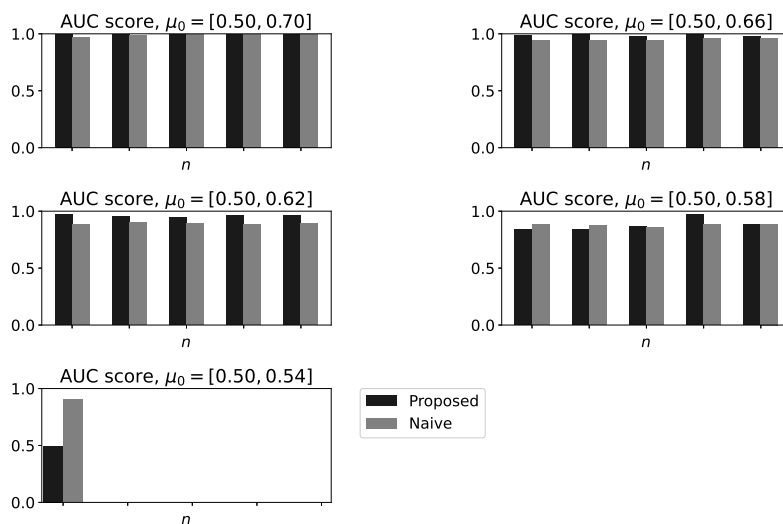


Figure 5.A9: AUC score in the realizable case for the **approximated** version of the proposed and naive statistical tests based on the biclustering algorithm using DeepTMR. If there were no null (i.e., $\hat{g} = g^{(N)}$) or alternative (i.e., $\hat{g} \neq g^{(N)}$) cases, respectively, then the corresponding bars would not have been plotted.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

15. Our dreams with my spouse are similar and harmonious.
16. We're compatible with my spouse about what love should be.
17. We share the same views about being happy in our life with my spouse.
18. My spouse and I have similar ideas about how marriage should be.
19. My spouse and I have similar ideas about how roles should be in marriage.
20. My spouse and I have similar values in trust.
21. I know exactly what my wife likes.
22. I know how my spouse wants to be taken care of when she/he sick.
23. I know my spouse's favorite food.
24. I can tell you what kind of stress my spouse is facing in her/his life.
25. I have knowledge of my spouse's inner world.
26. I know my spouse's basic anxieties.
27. I know what my spouse's current sources of stress are.
28. I know my spouse's hopes and wishes.
29. I know my spouse very well.
30. I know my spouse's friends and their social relationships.
31. I feel aggressive when I argue with my spouse.
32. When discussing with my spouse, I usually use expressions such as "you always" or "you never."
33. I can use negative statements about my spouse's personality during our discussions.
34. I can use offensive expressions during our discussions.
35. I can insult my spouse during our discussions.
36. I can be humiliating when we discussions.
37. My discussion with my spouse is not calm.
38. I hate my spouse's way of open a subject.

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

39. Our discussions often occur suddenly.
40. We're just starting a discussion before I know what's going on.
41. When I talk to my spouse about something, my calm suddenly breaks.
42. When I argue with my spouse, I only go out and I don't say a word.
43. I mostly stay silent to calm the environment a little bit.
44. Sometimes I think it's good for me to leave home for a while.
45. I'd rather stay silent than discuss with my spouse.
46. Even if I'm right in the discussion, I stay silent to hurt my spouse.
47. When I discuss with my spouse, I stay silent because I am afraid of not being able to control my anger.
48. I feel right in our discussions.
49. I have nothing to do with what I've been accused of.
50. I'm not actually the one who's guilty about what I'm accused of.
51. I'm not the one who's wrong about problems at home.
52. I wouldn't hesitate to tell my spouse about her/his inadequacy.
53. When I discuss, I remind my spouse of her/his inadequacy.
54. I'm not afraid to tell my spouse about her/his incompetence.

5.C Correspondence of the indices with locations in the metropolis traffic census dataset

The meaning of each row or column index of the metropolis traffic census dataset [1] is as follows:

- **[Tokyo]** 1: Chiyoda-ku, 2: Chuo-ku, 3: Minato-ku, 4: Shinjuku-ku, 5: Bunkyo-ku, 6: Taito-ku, 7: Sumida-ku, 8: Koto-ku, 9: Shinagawa-ku, 10: Meguro-ku, 11: Ota-ku, 12: Setagaya-ku, 13: Shibuya-ku, 14: Nakano-ku, 15: Suginami-ku, 16: Toshima-ku, 17: Kita-ku, 18: Arakawa-ku, 19: Itabashi-ku, 20: Nerima-ku, 21: Adachi-ku, 22: Katsushika-ku, 23: Edogawa-ku, 24: Hachioji City, 25: Tachikawa City, 26:

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Musashino City, 27: Mitaka City, 28: Ome City, 29: Fuchu City, 30: Akishima City, 31: Chofu City, 32: Machida City, 33: Koganei City, 34: Kodaira City, 35: Hino City, 36: Higashimurayama City, 37: Kokubunji City, 38: Kunitachi City, 39: Fussa City, 40: Komae City, 41: Higashiyamato City, 42: Kiyose City, 43: Higashikurume City, 44: Musashimurayama City, 45: Tama City, 46: Inagi City, 47: Hamura City, 48: Akiruno City, 49: Nishitokyo City, 50: Mizuho Town, 51: Hinode Town, 52: Hinohara Village, 53: Okutama Town

- **[Yokohama City, Kanagawa Prefecture]** 54: Tsurumi-ku, 55: Kanagawa-ku, 56: Nishi-ku, 57: Naka-ku, 58: Minami-ku, 59: Hodogaya-ku, 60: Isogo-ku, 61: Kanazawa-ku, 62: Kohoku-ku, 63: Totsuka-ku, 64: Konan-ku, 65: Asahi-ku, 66: Midori-ku, 67: Seya-ku, 68: Sakae-ku, 69: Izumi-ku, 70: Aoba-ku, 71: Tsuzuki-ku
- **[Kawasaki City, Kanagawa Prefecture]** 72: Kawasaki-ku, 73: Saiwai-ku, 74: Nakahara-ku, 75: Takatsu-ku, 76: Tama-ku, 77: Miyamae-ku, 78: Asao-ku
- **[Sagamihara City, Kanagawa Prefecture]** 79: Midori-ku, 80: Chuo-ku, 81: Minami-ku
- **[Kanagawa Prefecture]** 82: Yokosuka City, 83: Hiratsuka City, 84: Kamakura City, 85: Fujisawa City, 86: Odawara City, 87: Chigasaki City, 88: Zushi City, 89: Miura City, 90: Hadano City, 91: Atsugi City, 92: Yamato City, 93: Isehara City, 94: Ebina City, 95: Zama City, 96: Minamiashigara City, 97: Ayase City, 98: Hayama Town, 99: Samukawa Town, 100: Oiso Town, 101: Ninomiya Town, 102: Nakai Town, 103: Oimachi, 104: Matsuda Town, 105: Kaisei Town, 106: Hakone Town, 107: Manazuru Town, 108: Yugawara Town, 109: Aikawa Town
- **[Saitama City, Saitama Prefecture]** 110: Nishi-ku, 111: Kita-ku, 112: Omiya-ku, 113: Minuma-ku, 114: Chuo-ku, 115: Sakura-ku, 116: Urawa-ku, 117: Minami-ku, 118: Midori-ku, 119: Iwatsuki-ku
- **[Saitama Prefecture]** 120: Kawagoe City, 121: Kumagaya City, 122: Kawaguchi City, 123: Gyoda City, 124: Chichibu City, 125: Tokorozawa City, 126: Hanno City, 127: Kazo City, 128: Honjo City, 129: Higashimatsuyama City, 130: Kasukabe City, 131: Sayama City, 132: Hanyu City, 133: Konosu City, 134: Fukaya City, 135: Ageo City, 136: Soka City, 137: Koshigaya City, 138: Warabi City, 139: Toda City, 140: Iruma City, 141: Asaka City, 142: Shiki City, 143: Wako City, 144: Niiza City, 145: Okegawa City, 146: Kuki City, 147: Kitamoto City, 148: Yashio City, 149: Fujimi City, 150: Misato City, 151: Hasuda City, 152: Sakado City, 153: Satte City, 154: Tsurugashima City, 155: Hidaka City, 156: Yoshikawa City, 157: Fujimino City, 158: Shiraoka City, 159: Ina Town, 160: Miyoshi Town, 161: Moroyama Town, 162: Ogose Town, 163: Namegawa Town, 164: Ranzan Town, 165: Ogawa Town, 166: Kawajima Town, 167: Yoshimi Town, 168: Hatoyama Town, 169: Tokigawa

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

Town, 170: Yokoze Town, 171: Higashi Chichibu Village, 172: Misato Town, 173: Kamisato Town, 174: Yorii Town, 175: Miyashiro Town, 176: Sugito Town, 177: Matsubushi Town

- **[Chiba City, Chiba Prefecture]** 178: Chuo-ku, 179: Hanamigawa-ku, 180: Inage-ku, 181: Wakaba-ku, 182: Midori-ku, 183: Mihama-ku
- **[Chiba Prefecture]** 184: Ichikawa City, 185: Funabashi City, 186: Kisarazu City, 187: Matsudo City, 188: Noda City, 189: Mobara City, 190: Narita City, 191: Sakura City, 192: Togane City, 193: Narashino City, 194: Kashiwa City, 195: Ichihara City, 196: Nagareyama City, 197: Yachiyo City, 198: Abiko City, 199: Kamagaya City, 200: Kimitsu City, 201: Urayasu City, 202: Yotsukaido City, 203: Sodegaura City, 204: Yachimata City, 205: Inzai City, 206: Shiroy City, 207: Tomisato City, 208: Katori City, 209: Sanmu City, 210: Oamishirasato City, 211: Shisui Town, 212: Sakae Town, 213: Kozaki Town, 214: Ichinomiya Town, 215: Chosei Village, 216: Nagara Town, 217: Otaki Town
- **[Ibaraki Prefecture]** 218: Tsuchiura City, 219: Koga City, 220: Ishioka City, 221: Yuki City, 222: Ryugasaki City, 223: Shimotsuma City, 224: Joso City, 225: Toride City, 226: Ushiku City, 227: Tsukuba City, 228: Moriya City, 229: Chikusei City, 230: Bando City, 231: Inashiki City, 232: Kasumigaura City, 233: Tsukubamirai City, 234: Miho Village, 235: Ami Town, 236: Kawachi Town, 237: Yachiyo Town, 238: Goka Town, 239: Sakai Town, 240: Tone Town
- **[Gunma Prefecture]** 241: Tatebayashi City, 242: Itakura Town, 243: Meiwa Town
- **[Tochigi Prefecture]** 244: Tochigi City, 245: Sano City, 246: Oyama City, 247: Nogi Town
- **[Yamanashi Prefecture]** 248: Otsuki City, 249: Uenohara City

Conclusion

Summary and Follow-up Work

The summary for the three main contributions, which have been stated in Section 1.3, is as follows (see also Figure 1).

- **Evaluation of the number of biclusters:** In Chapter 3, we developed an asymptotic statistical test on the number of row and column clusters for an LBM. By sequentially testing the hypothetical cluster numbers in ascending order, we can estimate the number of biclusters in a given matrix. We can use an arbitrary biclustering algorithm for estimating the bicluster structure, as long as it satisfies the consistency condition (Section 3.2). The proposed test is asymptotically valid in the sense that the test statistic converges in law to the TW_1 distribution in the limit of matrix size $m \rightarrow \infty$. There is no theoretical guarantee of the proposed test for a finite size matrix, however, in Section 3.5.3, we demonstrated its effectiveness experimentally with the matrix sizes of several hundreds times several hundreds.
- **Evaluation of the estimated bicluster structure:** In Chapter 4, we proposed a selective test for the estimated bicluster structure of a given matrix based on the squared residue criterion. Unlike the asymptotic test in Chapter 3, we derived the exact null distribution with a finite size matrix. We can also test more detailed information (i.e., bicluster assignments) in the proposed method than in that of Chapter 3. However, to fulfill these conditions, stricter assumptions were required, including that each entry follows a Gaussian distribution and that the estimated bicluster structure is the global optimal solution of squared residue minimization problem. The latter condition makes the proposed test intractable with increasing matrix size m , since the number of possible bicluster assignments increases exponentially with m . To alleviate this problem, we also proposed an approximated test based on simulated annealing.
- **Extraction of the row and column features used for matrix reordering:** In Chapter 5, we constructed a new neural network model for matrix reordering. Unlike the existing methods, the proposed model extracts the row and column features

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

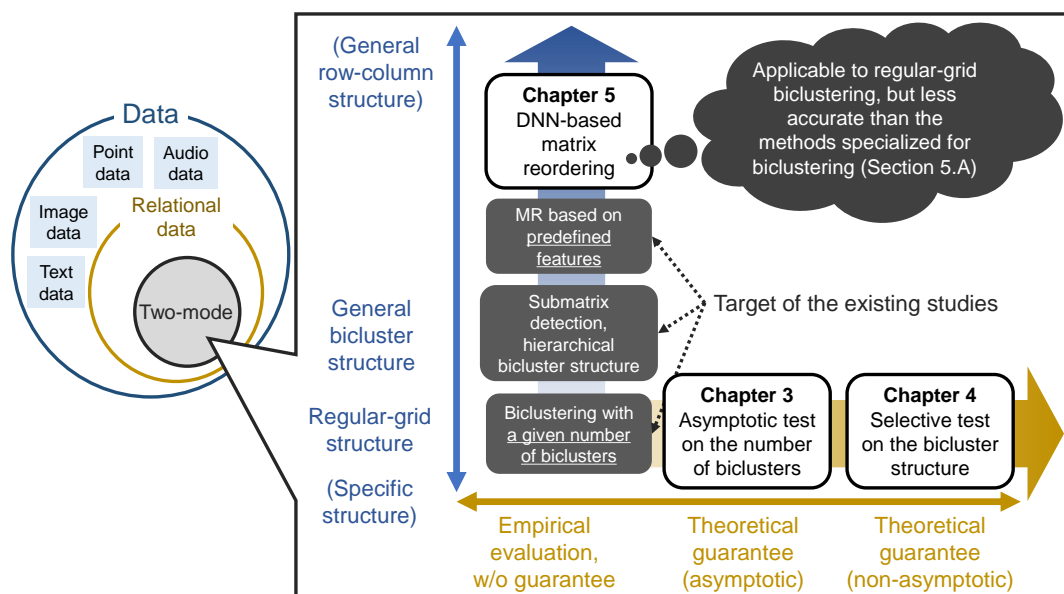


Figure 1: Summary of the main contributions of this dissertation.

in a data-driven way, based on which the row and column orderings are defined. The “goodness” of such row and column features is evaluated by reconstruction error of the model output, which can be seen as a denoised mean version of the observed matrix. On one hand, the proposed method can extract various types of latent structural patterns in a matrix without prior knowledge. On the other hand, it is affected the random effect in training (e.g., parameter initialization of the model).

Follow-up work As we have introduced in the last paragraph in Section 1.3, we have two follow-up studies for this dissertation. In [153], we develop a statistical test on the number of biclusters in more general problem setting than that of Chapter 3 in several ways. First, we consider the submatrix detection/localization problem in Section 2.2.2 instead of the regular-grid biclustering in Section 2.2.1. As we have discussed in Section 2.2.2, the former problem includes more general bicluster structures than the latter one. In the former setting, we assume that a given matrix consists of K biclusters and H background submatrices, where the entries in all the background submatrices follow identical distributions. Second, we assume that the number of biclusters (including background submatrices) might increase with matrix size m . Specifically, we show that the proposed test is asymptotically valid under the condition that $K + H = O\left(m^{\frac{1}{42}-\epsilon}\right)$ for some $\epsilon > 0$ in realizable case. Third, in [151], we propose a similar neural-network-based matrix reordering method for one-mode relational data (the definitions of two-mode and one-mode relational data are given in Sections 2.1.1 and 2.1.2, respectively). For one-mode relational data, we assume that

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

the orders of rows and columns are identical, regardless of the data matrix is symmetric or not (i.e., the network is directed or undirected). To fulfill this condition, we use a weight-sharing autoencoder network, which extracts the same node feature for the i th row and the i th column of a given adjacency matrix. These studies have extended this dissertation in some aspects, however, there still remains a room for improvement, as we describe in the next section.

Future Perspective

Aside from the limitation of each of the three contributions, which has been stated in Sections 3.6, 4.5, and 5.5, respectively, possible extensions for this dissertation are as follows.

- As for the statistical test on the number of biclusters, in Chapter 3, we derived the asymptotic behavior of the proposed test statistic T in both realizable and unrealizable cases. However, we have not considered the case in which a given matrix does not have any latent block structure (i.e., no block structure exists under which each entry independently follows a block-wise identical distribution). This case was also out of scope in the selective test of Chapter 4. When applying the proposed tests, it is considered as a special case of the alternative hypotheses, where we have no theoretical guarantee for now. It would be helpful to check the asymptotic behavior of the proposed test statistics in (at least a part of) this case through empirical simulation or to develop new test that can be applied in more general settings.
- The proposed statistical tests in Chapters 3 and 4 cannot be applied to sparse data, which contain unknown entries. Such a sparse data matrix can be transformed into a real matrix by substituting the unknown entries with some constant value (e.g., zero). However, since such substituted entries violate the LBM assumption that it is generated from a block-wise identical distribution, the (asymptotic) null distribution of the proposed test statistic is no longer valid. Future studies should analyze the null distribution of the proposed (or new) test statistic in this case.
- Another important future direction would be to provide some theoretical guarantee for the proposed matrix reordering method in Chapter 5, as well as the biclustering problem. There has been no statistical test directly on the row and column orderings. However, for the spectral and dimension-reduction methods, we can at least apply a statistical inference on the model of a given observed matrix (e.g., rank of the matrix) [2, 31, 81, 83]. It would be useful if we can develop a statistical inference method also for the proposed neural-network-based approach. Another theoretical

5. Matrix reordering method for capturing flexible structural patterns in a relational data matrix

direction would be to evaluate the denoising performance or reordering accuracy of the proposed DeepTMR model.

- In the original matrix reordering problem, we only consider a single set of row and column orderings that yields some latent pattern of a given matrix. It would be interesting to further extend this setting by combining the ideas of biclustering and matrix reordering, that is, by assuming the existence of multiple biclusters with different structural patterns. It would be useful if we can develop a statistical inference method for the number of biclusters or the bicluster structure for such a setting.
- Aside from biclustering and matrix reordering, we introduced various tasks with regard to a relational data matrix in Section 1.1. It would be further interesting if we can provide some meta-evaluation method for selecting an appropriate task (i.e., what type of information we should extract) for a given matrix.

Bibliography

- [1] e-Stat. https://www.e-stat.go.jp/stat-search/files?page=1&query=%E8%A1%8C%E6%94%BF%E5%8C%BA%E7%94%BB%E9%96%93%E7%A7%BB%E5%8B%95%E4%BA%BA%E5%93%A1%E8%A1%A8&layout=dataset&stat_infid=000031598030.
- [2] M. M. Al-Sadoon. A unifying theory of tests of rank. *Journal of Econometrics*, 199(1):49–62, 2017.
- [3] R. Albert, H. Jeong, and A.-L. Barabási. Diameter of the world-wide web. *Nature*, 401:130–131, 1999.
- [4] B. P. W. Ames. Guaranteed clustering and biclustering via semidefinite programming. *Mathematical Programming*, 147(1):429–465, 2014.
- [5] P. Arabie, S. A. Boorman, and P. R. Levitt. Constructing blockmodels: How and why. *Journal of Mathematical Psychology*, 17(1):21–63, 1978.
- [6] P. Arabie and L. J. Hubert. Combinatorial data analysis. *Annual Review of Psychology*, 43(1):169–203, 1992.
- [7] Z. D. Bai and Y. Q. Yin. Limit of the smallest eigenvalue of a large dimensional sample covariance matrix. *The Annals of Probability*, 21(3):1275–1294, 1993.
- [8] Z. Bao, G. Pan, and W. Zhou. Universality for the largest eigenvalue of sample covariance matrices with general population. *The Annals of Statistics*, 43(1):382–421, 2015.
- [9] Z. Bar-Joseph, D. K. Gifford, and T. S. Jaakkola. Fast optimal leaf ordering for hierarchical clustering. *Bioinformatics*, 17(suppl_1):S22–S29, 2001.
- [10] D. Basu. On statistics independent of a complete sufficient statistic. *Sankhyā: The Indian Journal of Statistics*, 15(4):377–380, 1955.

BIBLIOGRAPHY

- [11] M. Behrisch, B. Bach, N. H. Riche, T. Schreck, and J. Fekete. Matrix reordering methods for table and network visualization. *Computer Graphics Forum*, 35(3):693–716, 2016.
- [12] A. K. Bera and C. M. Jarque. Efficient tests for normality, homoscedasticity and serial independence of regression residuals: Monte Carlo Evidence. *Economics Letters*, 7(4):313–318, 1981.
- [13] R. Berk, L. Brown, A. Buja, K. Zhang, and L. Zhao. Valid post-selection inference. *The Annals of Statistics*, 41(2):802–837, 2013.
- [14] M. W. Berry, S. T. Dumais, and G. W. O’Brien. Using linear algebra for intelligent information retrieval. *SIAM Review*, 37(4):573–595, 1995.
- [15] M. W. Berry, B. Hendrickson, and P. Raghavan. Sparse matrix reordering schemes for browsing hypertext. *Lectures in Applied Mathematics*, 32(2):99–123, 1996.
- [16] P. J. Bickel and P. Sarkar. Hypothesis testing for automated community detection in networks. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(1):253–273, 2016.
- [17] A. Bloemendal, A. Knowles, H.-T. Yau, and J. Yin. On the principal components of sample covariance matrices. *Probability Theory and Related Fields*, 164:459–552, 2016.
- [18] A. Bloemendal, E. László, K. Antti, Y. Horng-Tzer, and Y. Jun. Isotropic local laws for sample covariance and generalized Wigner matrices. *Electronic Journal of Probability*, 19:1–53, 2014.
- [19] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008, 2008.
- [20] I. Borg and P. Groenen. *Modern Multidimensional Scaling: Theory and applications*. Springer Series in Statistics. Springer, New York, NY, 1997.
- [21] S. P. Borgatti. Centrality and network flow. *Social Networks*, 27(1):55–71, 2005.
- [22] K. O. Bowman and L. R. Shenton. Omnibus test contours for departures from normality based on $\sqrt{b_1}$ and b_2 . *Biometrika*, 62(2):243–250, 1975.
- [23] V. Brault and A. Channarond. Fast and consistent algorithm for the latent block model. arXiv:1610.09005, 2016.

BIBLIOGRAPHY

- [24] M. J. Brusco and S. Stahl. Using quadratic assignment methods to generate initial permutations for least-squares unidimensional scaling of symmetric proximity matrices. *Journal of Classification*, 17(2):197–223, 2000.
- [25] S. Busygin, O. Prokopyev, and P. M. Pardalos. Biclustering in data mining. *Computers & Operations Research*, 35(9):2964–2987, 2008.
- [26] C. Butucea and Y. I. Ingster. Detection of a sparse submatrix of a high-dimensional noisy matrix. *Bernoulli*, 19(5B):2652–2688, 2013.
- [27] T. T. Cai and Y. Wu. Statistical and computational limits for sparse matrix detection. *The Annals of Statistics*, 48(3):1593–1614, 2020.
- [28] G. Caraux and S. Pinloche. PermutMatrix: A graphical environment to arrange gene expression profiles in optimal linear order. *Bioinformatics*, 21(7):1280–1281, 2005.
- [29] B. Chen, F. Li, S. Chen, R. Hu, and L. Chen. Link prediction based on non-negative matrix factorization. *PLOS ONE*, 12(8):1–18, 2017.
- [30] K. Chen and J. Lei. Network cross-validation for determining the number of communities in network data. *Journal of the American Statistical Association*, 113(521):241–251, 2018.
- [31] Q. Chen and Z. Fang. Improved inference on the rank of a matrix. *Quantitative Economics*, 10(4):1787–1824, 2019.
- [32] Y. Cheng and G. M. Church. Biclustering of expression data. In *Proceedings. International Conference on Intelligent Systems for Molecular Biology*, volume 8, pages 93–103, 2000.
- [33] E. C. Chi, G. I. Allen, and R. G. Baraniuk. Convex biclustering. *Biometrics*, 73(1):10–19, 2017.
- [34] H. Cho, I. S. Dhillon, Y. Guan, and S. Sra. Minimum sum-squared residue co-clustering of gene expression data. In *Proceedings of the 2004 SIAM International Conference on Data Mining*, pages 114–125, 2004.
- [35] A. Clauset, M. E. J. Newman, and C. Moore. Finding community structure in very large networks. *Physical Review E*, 70(6):066111, 2004.
- [36] W. J. Conover. *Practical Nonparametric Statistics*. John Wiley & Sons, New York, 1999.

BIBLIOGRAPHY

- [37] M. Corneli, P. Latouche, and F. Rossi. Exact ICL maximization in a non-stationary time extension of the latent block model for dynamic networks. In *Proceedings of the 23-th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pages 225–230, 2015.
- [38] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, 1989.
- [39] B. Dabbs and B. Junker. Comparison of cross-validation methods for stochastic block models. arXiv:1605.03000, 2016.
- [40] G. Derval, V. Branders, P. Dupont, and P. Schaus. The maximum weighted submatrix coverage problem: A CP approach. In *Integration of Constraint Programming, Artificial Intelligence, and Operations Research*, pages 258–274, 2019.
- [41] I. S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 269–274, 2001.
- [42] J. Díaz, J. Petit, and M. Serna. A survey of graph layout problems. *ACM Computing Surveys*, 34(3):313–356, 2002.
- [43] X. Ding and F. Yang. A necessary and sufficient condition for edge universality at the largest singular values of covariance matrices. *The Annals of Probability*, 28(3):1679–1738, 2018.
- [44] D. Dua and C. Graff. UCI machine learning repository. <http://archive.ics.uci.edu/ml>, 2017. University of California, Irvine, School of Information and Computer Sciences.
- [45] M. B. Eisen, P. T. Spellman, P. O. Brown, and D. Botstein. Cluster analysis and display of genome-wide expression patterns. *Proceedings of the National Academy of Sciences*, 95(25):14863–14868, 1998.
- [46] W. Fithian, D. Sun, and J. Taylor. Optimal inference after model selection. arXiv:1410.2597, 2014.
- [47] C. J. Flynn and P. O. Perry. Profile likelihood biclustering. *Electronic Journal of Statistics*, 14(1):731–768, 2020.
- [48] S. Fortunato and M. Barthélemy. Resolution limit in community detection. *Proceedings of the National Academy of Sciences*, 104(1):36–41, 2007.
- [49] L. C. Freeman. A set of measures of centrality based on betweenness. *Sociometry*, 40(1):35–41, 1977.

BIBLIOGRAPHY

- [50] L. C. Freeman. Centrality in social networks conceptual clarification. *Social Networks*, 1(3):215–239, 1978.
- [51] M. Friendly. Corrgrams: Exploratory displays for correlation matrices. *The American Statistician*, 56(4):316–324, 2002.
- [52] K.-I. Funahashi. On the approximate realization of continuous mappings by neural networks. *Neural Networks*, 2(3):183–192, 1989.
- [53] A. Gangrade, P. Venkatesh, B. Nazer, and V. Saligrama. Efficient near-optimal testing of community changes in balanced stochastic block models. In *Advances in Neural Information Processing Systems 32*, pages 10364–10375, 2019.
- [54] S. Geman. A limit theorem for the norm of random matrices. *The Annals of Probability*, 8(2):252–261, 1980.
- [55] J. E. Gentle. *Matrix Algebra: Theory, Computations, and Applications in Statistics*. Springer, New York, NY, 2007.
- [56] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, 2010.
- [57] K. Goldberg, T. Roeder, D. Gupta, and C. Perkins. Eigentaste: A constant time collaborative filtering algorithm. *Information Retrieval*, 4(2):133–151, 2001.
- [58] G. Govaert and M. Nadif. Clustering with block mixture models. *Pattern Recognition*, 36:463–473, 2003.
- [59] R. Guimerà and L. A. N. Amaral. Modeling the world-wide airport network. *The European Physical Journal B*, 38(2):381–385, 2004.
- [60] R. Guimerà, S. Mossa, A. Turtleschi, and L. A. N. Amaral. The worldwide air transportation network: Anomalous centrality, community structure, and cities’ global roles. *Proceedings of the National Academy of Sciences*, 102(22):7794–7799, 2005.
- [61] B. Hajek. Cooling schedules for optimal annealing. *Mathematics of Operations Research*, 13(2):311–329, 1988.
- [62] F. M. Harper and J. A. Konstan. The MovieLens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems*, 5(4):1–19, 2015.

BIBLIOGRAPHY

- [63] J. A. Hartigan. Direct clustering of a data matrix. *Journal of the American Statistical Association*, 67(337):123–129, 1972.
- [64] R. Henriques and S. C. Madeira. B_{Sig}: evaluating the statistical significance of biclustering solutions. *Data Mining and Knowledge Discovery*, 32:124–161, 2018.
- [65] P. W. Holland, K. B. Laskey, and S. Leinhardt. Stochastic blockmodels: First steps. *Social Networks*, 5:109–137, 1983.
- [66] J. Hu, H. Qin, T. Yan, and Y. Zhao. Corrected Bayesian information criterion for stochastic block models. *Journal of the American Statistical Association*, 115(532):1771–1783, 2020.
- [67] J. Hu, J. Zhang, H. Qin, T. Yan, and J. Zhu. Using maximum entry-wise deviation to test the goodness of fit for stochastic block models. *Journal of the American Statistical Association*, 116(535):1373–1382, 2021.
- [68] P. Ihm. A contribution to the history of seriation in archaeology. In *Classification — the Ubiquitous Challenge*, pages 307–316, 2005.
- [69] S. Inoue, Y. Umezū, S. Tsubota, and I. Takeuchi. Post clustering inference for heterogeneous data. In *Information-Based Induction Science Workshop*, pages 69–76, 2017.
- [70] B. Irie and S. Miyake. Capabilities of three-layered perceptrons. In *IEEE 1988 International Conference on Neural Networks*, volume 1, pages 641–648, 1988.
- [71] C. M. Jarque and A. K. Bera. Efficient tests for normality, homoscedasticity and serial independence of regression residuals. *Economics Letters*, 6(3):255–259, 1980.
- [72] K. Johansson. Shape fluctuations and random matrices. *Communications in Mathematical Physics*, 209:437–476, 2000.
- [73] I. M. Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *The Annals of Statistics*, 29(2):295–327, 2001.
- [74] V. Karwa, D. Pati, S. Petrović, L. Solus, N. Alexeev, M. Raič, D. Wilburne, R. Williams, and B. Yan. Exact tests for stochastic block models. arXiv:1612.06040, 2016.
- [75] T. Kawamoto and Y. Kabashima. Cross-validation estimate of the number of clusters in a network. *Scientific Reports*, 7(1):3327, 2017.
- [76] C. Keribin, V. Brault, G. Celeux, and G. Govaert. Model selection for the binary latent block model. In *Proceedings of 20th International Conference on Computational Statistics*, pages 379–390, 2012.

BIBLIOGRAPHY

- [77] C. Keribin, V. Brault, G. Celeux, and G. Govaert. Estimation and selection for the latent block model on categorical data. *Statistics and Computing*, 25:1201–1216, 2015.
- [78] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations*, 2015.
- [79] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220:671–680, 1983.
- [80] M. Kitsak, L. K. Gallos, S. Havlin, F. Liljeros, L. Muchnik, H. E. Stanley, and H. A. Makse. Identification of influential spreaders in complex networks. *Scientific Reports*, 6(11):888–893, 2010.
- [81] F. Kleibergen and R. Paap. Generalized reduced rank tests using the singular value decomposition. *Journal of Econometrics*, 133(1):97–126, 2006.
- [82] Y. Kluger, R. Basri, J. T. Chang, and M. Gerstein. Spectral biclustering of microarray data: Coclustering genes and conditions. *Genome Research*, 13(4):703–716, 2003.
- [83] K. Konstantinides and K. Yao. Statistical analysis of effective singular values in matrix rank determination. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 36(5):757–763, 1988.
- [84] L. Labiod and M. Nadif. Modularity and spectral co-clustering for categorical data. In *Proceedings of the International Conference on Data Mining*, pages 386–392, 2011.
- [85] D. D. Lee and H. S. Seung. Unsupervised learning by convex and conic coding. In *Advances in Neural Information Processing Systems 9*, pages 515–521, 1996.
- [86] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.
- [87] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in Neural Information Processing Systems 13*, pages 535–541, 2000.
- [88] J. D. Lee, D. L. Sun, Y. Sun, and J. E. Taylor. Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927, 2016.
- [89] J. D. Lee, Y. Sun, and J. E. Taylor. Evaluating the statistical significance of biclusters. In *Advances in Neural Information Processing Systems 28*, pages 1324–1332, 2015.
- [90] J. D. Lee and J. E. Taylor. Exact post model selection inference for marginal screening. In *Advances in Neural Information Processing Systems 27*, pages 136–144, 2014.

BIBLIOGRAPHY

- [91] M. Lee, H. Shen, J. Z. Huang, and J. S. Marron. Biclustering via sparse singular value decomposition. *Biometrics*, 66(4):1087–1095, 2010.
- [92] J. Lei. A goodness-of-fit test for stochastic block models. *The Annals of Statistics*, 44(1):401–424, 2016.
- [93] J. Leskovec, J. Kleinberg, and C. Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Transactions on Knowledge Discovery from Data*, 1(1):2–es, 2007.
- [94] J. Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney. Community structure in large networks: Natural cluster sizes and the absence of large well-defined clusters. *Internet Mathematics*, 6(1):29–123, 2009.
- [95] J. Leskovec and J. McAuley. Learning to discover social circles in ego networks. In *Advances in Neural Information Processing Systems 25*, pages 539–547, 2012.
- [96] J. Y.-T. Leung, O. Vornberger, and J. D. Witthoff. On some variants of the bandwidth minimization problem. *SIAM Journal on Computing*, 13(3):650–667, 1984.
- [97] T. Li, E. Levina, and J. Zhu. Network cross-validation by edge sampling. *Biometrika*, 107(2):257–276, 2020.
- [98] I. Liiv. Seriation and matrix reordering methods: An historical overview. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 3(2):70–91, 2010.
- [99] Y. Lin and J. Yuan. Profile minimization problem for matrices and graphs. *Acta Mathematicae Applicatae Sinica*, 10:107–112, 1994.
- [100] L. Liu, D. M. Hawkins, S. Ghosh, and S. S. Young. Robust singular value decomposition analysis of microarray data. *Proceedings of the National Academy of Sciences*, 100(23):13167–13172, 2003.
- [101] Y. Liu and J. Guo. Distribution-free, size adaptive submatrix detection with acceleration. arXiv:1804.10887, 2018.
- [102] Y. Liu, M. Tang, T. Zhou, and Y. Do. Core-like groups result in invalidation of identifying super-spreader by k-shell decomposition. *Scientific Reports*, 5(1):9602, 2015.
- [103] J. R. Loftus and J. E. Taylor. Selective inference in regression models with groups of variables. arXiv:1511.01478, 2015.
- [104] A. Lomet, G. Govaert, and Y. Grandvalet. Model selection in block clustering by the integrated classification likelihood. In *Proceedings of the 20th International Conference on Computational Statistics*, pages 519–530, 2012.

BIBLIOGRAPHY

- [105] Z. Ma. Accuracy of the Tracy-Widom limits for the extreme eigenvalues in white Wishart matrices. *Bernoulli*, 18(1):322–359, 2012.
- [106] Z. Ma and Y. Wu. Computational barriers in minimax submatrix detection. *The Annals of Statistics*, 43(3):1089–1116, 2015.
- [107] S. C. Madeira and A. L. Oliveira. Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 1(1):24–45, 2004.
- [108] M. Mariadassou and C. Matias. Convergence of the groups posterior distribution in latent or stochastic block models. *Bernoulli*, 21(1):537–573, 2015.
- [109] A. K. Menon and C. Elkan. Link prediction via matrix factorization. In *ECML PKDD 2011: Machine Learning and Knowledge Discovery in Databases*, pages 437–452, 2011.
- [110] F. Murtagh. Book review: W. Gaul and M. Schader, Eds., Data, expert knowledge and decisions, Heidelberg: Springer-Verlag, 1988, pp. viii + 380. *Journal of Classification*, 6:129–132, 1989.
- [111] M. Nadif and G. Govaert. Model-based co-clustering for continuous data. In *Proceedings of the 9th International Conference on Machine Learning and Applications*, pages 175–180, 2010.
- [112] M. Nakano, K. Ishiguro, A. Kimura, T. Yamada, and N. Ueda. Rectangular tiling process. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 361–369, 2014.
- [113] M. E. J. Newman. The structure of scientific collaboration networks. *Proceedings of the National Academy of Sciences*, 98(2):404–409, 2001.
- [114] M. E. J. Newman. Finding community structure in networks using the eigenvectors of matrices. *Physical Review E*, 74(3):036104, 2006.
- [115] M. E. J. Newman. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582, 2006.
- [116] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical Review E*, 69(2):026113, 2004.
- [117] F. S. Passino and N. A. Heard. Bayesian estimation of the latent dimension and communities in stochastic blockmodels. *Statistics and Computing*, 30:1291–1307, 2020.

BIBLIOGRAPHY

- [118] S. Péché. Universality results for the largest eigenvalues of some sample covariance matrix ensembles. *Probability Theory and Related Fields*, 143:481–516, 2009.
- [119] T. P. Peixoto. Parsimonious module inference in large networks. *Physical Review Letters*, 110:148701, 2013.
- [120] V. Perrone, P. A. Jenkins, D. Spanò, and Y. W. Teh. Poisson random fields for dynamic feature models. *Journal of Machine Learning Research*, 18(127):1–45, 2017.
- [121] W. M. F. Petrie. Sequences in prehistoric remains. *The Journal of the Anthropological Institute of Great Britain and Ireland*, 29(3/4):295–301, 1899.
- [122] N. S. Pillai and J. Yin. Universality of covariance matrices. *The Annals of Applied Probability*, 24(3):935–1001, 2014.
- [123] B. Pontes, R. Giráldez, and J. S. Aguilar-Ruiz. Biclustering on expression data: A review. *Journal of Biomedical Informatics*, 57:163–180, 2015.
- [124] R. Rastelli and N. Friel. Optimal Bayesian estimators for latent variable cluster models. *Statistics and Computing*, 28:1169–1186, 2018.
- [125] V. Robert, Y. Vasseur, and V. Brault. Comparing high-dimensional partitions with the co-clustering Adjusted Rand Index. *Journal of Classification*, 38:158–186, 2021.
- [126] W. S. Robinson. A method for chronologically ordering archaeological deposits. *American Antiquity*, 16(4):293–301, 1951.
- [127] J. L. Rodgers and T. D. Thompson. Seriation and multidimensional scaling: A data analysis approach to scaling asymmetric proximity matrices. *Applied Psychological Measurement*, 16(2):105–117, 1992.
- [128] D. M. Roy, C. Kemp, V. Mansinghka, and J. Tenenbaum. Learning annotated hierarchies from relational data. In *Advances in Neural Information Processing Systems 19*, pages 1185–1192, 2006.
- [129] D. M. Roy and Y. W. Teh. The Mondrian process. In *Advances in Neural Information Processing Systems 21*, pages 1377–1384, 2008.
- [130] D. F. S., Y. Yu, and Y. Feng. How many communities are there? *Journal of Computational and Graphical Statistics*, 26(1):171–181, 2017.
- [131] H. B. Saber, M. Elloumi, and M. Nadif. Block mixture model for the biclustering of microarray data. In *Proceedings of the 22nd International Workshop on Database and Expert Systems Applications*, pages 423–427, 2011.

BIBLIOGRAPHY

- [132] Y. Sakai and K. Yamanishi. An NML-based model selection criterion for general relational data modeling. In *2013 IEEE International Conference on Big Data*, pages 421–429, 2013.
- [133] B. M. Sarwar, G. Karypis, J. A. Konstan, and J. T. Riedl. Application of dimensionality reduction in recommender system - a case study. Technical report, Department of Computer Science and Engineering, University of Minnesota, 2000.
- [134] A. A. Shabalin, V. J. Weigman, C. M. Perou, and A. B. Nobel. Finding large average submatrices in high dimensional data. *The Annals of Applied Statistics*, 3(3):985–1012, 2009.
- [135] H. Shan and A. Banerjee. Bayesian co-clustering. In *Proceedings of the 8th IEEE International Conference on Data Mining*, pages 530–539, 2008.
- [136] J. Shao. *Mathematical Statistics*. Springer-Verlag New York, 2003.
- [137] J. W. Silverstein. The smallest eigenvalue of a large dimensional Wishart matrix. *The Annals of Probability*, 13(4):1364–1368, 1985.
- [138] A. Soshnikov. A note on universality of the distribution of the largest eigenvalues in certain sample covariance matrices. *Journal of Statistical Physics*, 108:1033–1056, 2002.
- [139] I. Spence and J. Graef. The determination of the underlying dimensionality of an empirically obtained matrix of proximities. *Multivariate Behavioral Research*, 9(3):331–341, 1974.
- [140] K. M. Tan and D. M. Witten. Sparse biclustering of transposable data. *Journal of Computational and Graphical Statistics*, 23:985–1008, 2014.
- [141] A. Tanay, R. Sharan, and R. Shamir. Biclustering algorithms: A survey. *Handbook of Computational Molecular Biology*, 9:1–20, 2005.
- [142] Y. Terada and H. Shimodaira. Selective inference for the problem of regions via multiscale bootstrap. arXiv:1711.00949, 2017.
- [143] X. Tian and J. Taylor. Selective inference with a randomized response. *The Annals of Statistics*, 46(2):679–710, 2018.
- [144] R. J. Tibshirani, A. Rinaldo, R. Tibshirani, and L. Wasserman. Uniform asymptotic inference and the bootstrap after model selection. *The Annals of Statistics*, 46(3):1255–1287, 2018.

BIBLIOGRAPHY

- [145] C. A. Tracy and H. Widom. The distributions of random matrix theory and their applications. In *New Trends in Mathematical Physics*, pages 753–765. Springer, 2009.
- [146] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- [147] V. Černý. Thermodynamical approach to the traveling salesman problem: An efficient simulation algorithm. *Journal of Optimization Theory and Applications*, 45:41–51, 1985.
- [148] J. D. Vlok and J. C. Olivier. Analytic approximation to the largest eigenvalue distribution of a white Wishart matrix. *IET Communications*, 6(12):1804–1811, 2012.
- [149] W. Wang, F. Cai, P. Jiao, and L. Pan. A perturbation-based framework for link prediction via non-negative matrix factorization. *Scientific Reports*, 6(1):38938, 2016.
- [150] J. H. Ward, Jr. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, 58(301):236–244, 1963.
- [151] C. Watanabe and T. Suzuki. AutoLL: Automatic linear layout of graphs based on deep neural network. In *IEEE Symposium Series on Computational Intelligence*, 2021.
- [152] C. Watanabe and T. Suzuki. Goodness-of-fit test for latent block models. *Computational Statistics & Data Analysis*, 154:107090, 2021. doi:10.1016/j.csda.2020.107090.
- [153] C. Watanabe and T. Suzuki. A goodness-of-fit test on the number of biclusters in a relational data matrix. arXiv:2102.11658, 2021.
- [154] C. Watanabe and T. Suzuki. Selective inference for latent block models. *Electronic Journal of Statistics*, 15(1):3137–3183, 2021. doi:10.1214/21-EJS1853.
- [155] C. Watanabe and T. Suzuki. Deep two-way matrix reordering for relational data analysis. *Neural Networks*, 146:303–315, 2022. doi:10.1016/j.neunet.2021.11.028.
- [156] D. J. Watts and S. H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393:440–442, 1998.
- [157] J. G. White, E. Southgate, J. N. Thomson, and S. Brenner. The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 314:1–340, 1986.

BIBLIOGRAPHY

- [158] J. Wyse and N. Friel. Block clustering with collapsed latent block models. *Statistics and Computing*, 22:415–428, 2012.
- [159] J. Wyse, N. Friel, and P. Latouche. Inferring structure in bipartite networks using the latent blockmodel and exact ICL. *Network Science*, 5(1):45–69, 2017.
- [160] K. Yamanishi, T. Wu, S. Sugawara, and M. Okada. The decomposed normalized maximum likelihood code-length criterion for selecting hierarchical latent variable models. *Data Mining and Knowledge Discovery*, 33(4):1017–1058, 2019.
- [161] Y. Q. Yin, Z. D. Bai, and P. R. Krishnaiah. On the limit of the largest eigenvalue of the large dimensional sample covariance matrix. *Probability Theory and Related Fields*, 78:509–521, 1988.
- [162] M. K. Yöntem. *The predictive role of the styles of parenthood origin on divorce predictors*. PhD thesis, Gaziosmanpasa University, 2017.
- [163] M. K. Yöntem, K. Adem, T. İlhan, and S. Kılıçarslan. Divorce prediction using correlation based feature selection and artificial neural networks. *Nevşehir HacıBektaş Veli Üniversitesi SBE Dergisi*, 9:259–273, 2019.
- [164] M. Yuan, Y. Feng, and Z. Shang. A likelihood-ratio type test for stochastic block models with bounded degrees. arXiv:1807.04426, 2018.
- [165] W. W. Zachary. An information flow model for conflict and fission in small groups. *Journal of Anthropological Research*, 33(4):452–473, 1977.
- [166] K. Zhao, Y. Rong, J. X. Yu, W. Huang, J. Huang, and H. Zhang. Graph ordering: Towards the optimal by learning. In *International Conference on Web Information Systems Engineering*, pages 423–437, 2021.
- [167] T. Zhou, J.-G. Liu, W.-J. Bai, G. Chen, and B.-H. Wang. Behaviors of susceptible-infected epidemics on scale-free networks with identical infectivity. *Physical Review E*, 74(5):056109, 2006.
- [168] H. Zhu, G. Mateos, G. B. Giannakis, N. D. Sidiropoulos, and A. Banerjee. Sparsity-cognizant overlapping co-clustering for behavior inference in social networks. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 3534–3537, 2010.