

Determining the Number of Factors in High-dimensional Generalised Latent Factor Models

Yunxiao Chen,

London School of Economics and Political Science

Xiaoou Li,

University of Minnesota

Abstract

As a generalisation of the classical linear factor model, generalised latent factor models are a useful tool for analysing multivariate data of different types, including binary choices and counts. In this paper, we propose an information criterion to determine the number of factors in generalised latent factor models. The consistency of the proposed information criterion is established under a high-dimensional setting where both the sample size and the number of manifest variables grow to infinity and data may have many missing values. To establish this consistency result, an error bound is established for the parameter estimates that improves the existing results and may be of independent theoretical interest. Simulation shows that the proposed criterion has good finite sample performance. An application to Eysenck's personality questionnaire confirms the three-factor structure of this personality survey.

KEYWORDS: Generalised latent factor model; Joint maximum likelihood estimator; High-dimensional data; Information criteria; Selection consistency

1 Introduction

Factor analysis is a popular method in social and behavioral sciences, including psychology, economics, and marketing (Bartholomew et al., 2011). It uses a relatively small number of factors to model the variation in a larger number of manifest variables. For example, in psychological science, manifest variables may correspond to personality questionnaire items for which factors are often interpreted as personality traits. Multivariate data in social and behavioral sciences often involve categorical or count variables, for which the classical linear factor model may not be suitable. Generalised latent factor models (Skron dal and Rabe-Hesketh, 2004; Chen et al., 2019b) provide a flexible framework for more types of data by combining generalised linear models and factor analysis. Specifically, item response theory models (Embretson and Reise, 2000; Reckase, 2009), which are widely used in psychological measurement and educational testing, can be viewed as special cases of generalised latent factor models.

Factor analysis is often used in an exploratory manner for generating scientific hypotheses, known as exploratory factor analysis. In that case, the number of factors and the corresponding loading structure are unknown and need to be learned from data. Quite a few methods have been proposed for determining the number of factors in linear factor models, including eigenvalue thresholding (Kaiser, 1960), subjective search for eigengap based on scree plot (Cattell, 1966), information criteria (Bai and Ng, 2002), cross-validation (Owen and Wang, 2016), and parallel analysis (Horn, 1965; Buja and Eyuboglu, 1992; Dobriban and Owen, 2019). However, fewer methods are available for determining the number of factors in generalised latent factor models and statistical theory remains to be developed, especially under a high-dimensional setting when the sample size and the number of manifest variables are large.

Traditionally, statistical inference of generalised latent factor models is usually carried out based on a marginal likelihood function (Bock and Aitkin, 1981; Skron dal and Rabe-Hesketh, 2004), in which latent factors are treated as random variables and integrated out. However,

for high-dimensional data involving large numbers of observations, manifest variables and factors, marginal-likelihood-based inference tends to suffer from a high computational burden and thus may not be always feasible. In that case, a joint likelihood function may be a good alternative that treats factors as fixed model parameters (Chen et al., 2019a,b; Zhu et al., 2016). Specifically, a joint maximum likelihood estimator is proposed in Chen et al. (2019a,b) that is not only easy to compute, but also statistically optimal in the minimax sense when both the sample size and the number of manifest variables grow to infinity. With a diverging number of parameters in the joint likelihood function, the classical information criteria, such as the Akaike information criterion (AIC; Akaike, 1974) and the Bayesian information criterion (BIC; Schwarz, 1978), may no longer be suitable.

In this paper, we propose a joint-likelihood-based information criterion (JIC) for determining the number of factors in generalised latent factor models. The proposed criterion is suitable for high-dimensional data with large numbers of observations and manifest variables, and can be used even when data contain many missing values. Under a very general setting, we prove the consistency of the proposed JIC when both the numbers of samples and manifest variables grow to infinity. Specifically, the missing entries are allowed to be non-uniformly distributed in the data matrix and their proportion is allowed to grow to one (i.e., the proportion of observable entries decays to zero). An error bound for the joint maximum likelihood estimator is established under a general missing data setting that improves upon the ones given in Cai and Zhou (2013), Davenport et al. (2014), Chatterjee et al. (2015), Bhaskar and Javanmard (2015), Ni and Gu (2016), and Chen et al. (2019b). Simulation shows that the proposed JIC has good finite sample performance under different settings and an application to the revised Eysenck’s personality questionnaire (Eysenck et al., 1985) finds three factors which confirms the design of this personality survey.

2 Joint-likelihood-based Information Criterion

2.1 Generalised Latent Factor Models

We consider multivariate data involving N individuals and J manifest variables. Let y_{ij} be a random variable denoting the i th individual's value on the j th manifest variable. Factor models assume that each individual is associated with K latent factors, denoted by a vector $F_i = (f_{i1}, \dots, f_{iK})^T$. We assume that the distribution of y_{ij} given F_i follows an exponential family distribution with natural parameter $d_j + A_j^T F_i$, and possibly a scale (i.e. dispersion) parameter ϕ , where d_j and $A_j = (a_{j1}, \dots, a_{jK})^T$ are manifest-variable-specific parameters. Specifically, d_j can be viewed as an intercept parameter, and a_{jk} is known as a loading parameter. More precisely, the probability density/mass function for y_{ij} takes the form

$$g(y|A_j, d_j, F_i, \phi) = \exp \left(\frac{y(d_j + A_j^T F_i) - b(d_j + A_j^T F_i)}{\phi} + c(y, \phi) \right), \quad (1)$$

where b and c are pre-specified functions that depend on the exponential family distribution. Given all the person- and manifest-variable-specific parameters, data y_{ij} , $i = 1, \dots, N$, $j = 1, \dots, J$, are assumed to be independent. In particular, linear factor models for continuous data, logistic for binary data, and Poisson model for counts, are special cases of model (1). We present the logistic and Poisson models as two examples, while pointing out that (1) also includes linear factor models as a special case when the exponential family distribution is chosen to be a Gaussian distribution.

Example 1. *When data are binary, (1) leads to a logistic model. That is, by letting $b(d_j + A_j^T F_i) = \log(1 + \exp(d_j + A_j^T F_i))$, $\phi = 1$, and $c(y, \phi) = 0$, (1) implies that y_{ij} follows Bernoulli distribution, with success probability $\exp(d_j + A_j^T F_i) / (1 + \exp(d_j + A_j^T F_i))$. This model is known as the multi-dimensional two-parameter logistic model (Reckase, 2009) that is widely used in educational testing and psychological measurement.*

Example 2. *For count data, (1) leads to a Poisson model by letting $b(d_j + A_j^T F_i) = \exp(d_j + A_j^T F_i)$*

$A_j^T F_i$, $\phi = 1$, and $c(y, \phi) = -\log(y!)$. Then y_{ij} follows a Poisson distribution with intensity $\exp(d_j + A_j^T F_i)$. This model is known as the Poisson factor model for count data (Wedel et al., 2003).

We further take missing data into account, under an ignorable missingness assumption. Let ω_{ij} be a binary random variable, indicating the missingness of y_{ij} . Specifically, $\omega_{ij} = 1$ means that y_{ij} is observed and $\omega_{ij} = 0$ if y_{ij} is missing. It is assumed that, given all the person- and manifest-variable-specific parameters, the missing indicators ω_{ij} , $i = 1, \dots, N, j = 1, \dots, J$, are independent of each other, and are also independent of data y_{ij} . The same missing data setting is adopted in Cai and Zhou (2013) for a 1-bit matrix completion problem and Zhu et al. (2016) for collaborative filtering. For nonignorable missing data, one may need to model the distribution of ω_{ij} given y_{ij} , F_i , A_j , and d_j . See Little and Rubin (2019) for more discussions on non-ignorable missingness. For the ease of explanation, in what follows, we assume the dispersion parameter $\phi > 0$ is known and does not change with N and J . Our theoretical development below can be extended to the case when ϕ is unknown; see Remark 4 below for a discussion.

2.2 Proposed Information Criterion

Under the above setting for generalised latent factor models, the log-likelihood function for observed data takes the form

$$l_K(F_1, \dots, F_N, A_1, d_1, \dots, A_J, d_J) = \sum_{\omega_{ij}=1} \log g(y_{ij} | A_j, d_j, F_i, \phi).$$

Note that a subscript K is added to the likelihood function to emphasize the number of factors in the current model.

For exploratory factor analysis, we consider the following constrained joint maximum

likelihood estimator as proposed in Chen et al. (2019a,b)

$$\begin{aligned}
(\hat{F}_1, \dots, \hat{F}_N, \hat{A}_1, \hat{d}_1, \dots, \hat{A}_J, \hat{d}_J) \in & \arg \max l_K(F_1, \dots, F_N, A_1, d_1, \dots, A_J, d_J), \\
s.t. & (\|F_i\|^2 + 1)^{\frac{1}{2}} \leq C, i = 1, \dots, N, \\
& (d_j^2 + \|A_j\|^2)^{\frac{1}{2}} \leq C, j = 1, \dots, J,
\end{aligned} \tag{2}$$

where $\|\cdot\|$ denotes the standard Euclidian norm. Here C is a reasonably large constant to ensure that a finite solution to (2) exists and satisfies certain regularity conditions.

As there is no further constraint imposed under the exploratory factor analysis setting, the solution to (2) is not unique. This indeterminacy of the solution will not be an issue for determining the number of factors, since the proposed JIC only depends on the log-likelihood function value rather than specific parameter values. The computation of (2) can be done by an alternating maximisation algorithm which has good convergence properties according to numerical experiments (Chen et al., 2019a,b), even though (2) is a non-convex optimization problem.

Let n be the number of observed data entries, i.e.,

$$n = \sum_{i=1}^N \sum_{j=1}^J \omega_{ij}.$$

The proposed JIC takes the form

$$\text{JIC}(K) = -2l_K(\hat{F}_1, \dots, \hat{F}_N, \hat{A}_1, \hat{d}_1, \dots, \hat{A}_J, \hat{d}_J) + v(n, N, J, K),$$

where \hat{F}_i , \hat{A}_j , and \hat{d}_j are given by (2), and $v(n, N, J, K)$ is a penalty term depending on n , N , J , and K that will be discussed below. We choose \hat{K} that minimizes $\text{JIC}(K)$.

As will be shown in Section 3, the consistency of \hat{K} can be guaranteed under a wide

range choice of $v(n, N, J, K)$. In practice, we suggest to use

$$v(n, N, J, K) = K(N \vee J) \log(n/(N \vee J)), \quad (3)$$

where $N \vee J$ denotes the maximum of N and J . When there is no missing data, i.e., $n = NJ$, then (3) becomes $v(n, N, J, K) = K(N \vee J) \log(N \wedge J)$, where $N \wedge J$ denotes the minimum of N and J . The advantage of this choice will be clarified in Section 3.

3 Theoretical Results

We start with the definition of several useful quantities. Let $p_{ij} = \Pr(\omega_{ij} = 1)$ be the sampling weight for y_{ij} and $p_{\min} = \min_{1 \leq i \leq N, 1 \leq j \leq J} p_{ij}$ be their minimum. Also let $n^* = \sum_{i=1}^N \sum_{j=1}^J p_{ij}$, $n_{i\cdot}^* = \sum_{j=1}^J p_{ij}$, and $n_{\cdot j}^* = \sum_{i=1}^N p_{ij}$ be the expected number of observations in the entire data matrix, each row and each column, respectively. Let $p_{\max} = (J^{-1} \max_{1 \leq i \leq N} n_{i\cdot}^*) \vee (N^{-1} \max_{1 \leq j \leq J} n_{\cdot j}^*)$ be the maximum average sampling weights for different columns and rows. Let $m_{ij}^* = d_j^* + (A_j^*)^T F_i^*$ be the true natural parameter for y_{ij} , and let $M^* = (m_{ij}^*)_{1 \leq i \leq N, 1 \leq j \leq J}$. We also denote $\hat{M} = (\hat{d}_j + \hat{A}_j^T \hat{F}_i)_{N \times J}$ to be the corresponding estimator of M obtained from (2). To emphasize the dependence on the number of factors, we use $\hat{M}^{(K)}$ to denote the estimator when assuming K factors in the model. Let K_{\max} denote the maximum number of factors considered in the model selection process and let K^* be the true number of factors.

The following two assumptions are made throughout the paper.

Assumption 1. For all $x \in [-2C^2, 2C^2]$, $b(x) < \infty$.

Assumption 2. The true model parameters F_i^* , A_j^* , and d_j^* satisfy the constraint in (2). That is, $(\|F_i^*\|^2 + 1)^{\frac{1}{2}} \leq C$ and $((d_j^*)^2 + \|A_j^*\|^2)^{\frac{1}{2}} \leq C$, for all i and j .

In the rest of the section, we will first present error bounds for the joint maximum likelihood estimator, and then present conditions on $v(n, N, J, K)$ that guarantee consistent

model selection.

Theorem 1. Assume $n^*/(\log n^*)^2 \geq (N \wedge J) \log(N + J)$ and the true number of factors satisfies $1 \leq K^* \leq K_{\max}$. Then, there is a finite constant κ depending on p_{\max}/p_{\min} , C , ϕ , the function b and independent with K_{\max} , N , J and n^* , such that with probability $1 - (n^*)^{-1} - 2(N + J)^{-1}$,

$$\max_{K^* \leq K \leq K_{\max}} \left\{ (NJ)^{-1/2} \|\hat{M}^{(K)} - M^*\|_F \right\} \leq \kappa \left\{ \frac{K_{\max}(N \vee J)}{n^*} \right\}^{1/2}. \quad (4)$$

In particular, if K^* is known then we have $(NJ)^{-1/2} \|\hat{M}^{(K^*)} - M\|_F \leq \kappa \left\{ \frac{K^*(N \vee J)}{n^*} \right\}^{1/2}$.

We make a few remarks on the above theorem.

Remark 1. It is well-known that in exploratory factor analysis, the factors F_1, \dots, F_N are not identifiable due to rotational indeterminacy, while m_{ij} s are identifiable. Thus, we establish error bounds for estimating the matrix M rather than those of F_i s and A_j s. If additional design information is available and a confirmatory generalised latent factor model is used, then the methods described in Section 2.2 and theoretical results in Theorem 1 may be extended to establish error bounds for F_i s following a similar strategy as in Chen et al. (2019b).

Remark 2. The error bound (4) improves several recent results in low-rank matrix estimation and completion. For example, when $n^* = o((N \wedge J)^2)$, it improves the error rate $O(\{(N \vee J)(n^*)^{-1} + NJ(n^*)^{-3/2}\}^{1/2})$ in Chen et al. (2019b) where fixed K^* and uniform sampling (i.e., $p_{\max} = p_{\min}$) are assumed. Other examples include Bhaskar and Javanmard (2015) and Ni and Gu (2016) where the error rates are shown to be $O(\{K^*(N \vee J) \log(N + J)(n^*)^{-1}\}^{1/2})$ and $O(K^*(N \vee J)^{1/2}(n^*)^{-1/2} + (N \vee J)^3(N \wedge J)^{1/2}(K^*)^{3/2}(n^*)^{(-2)})$ assuming binary data. The error estimate (4) is also smaller than the optimal rate $\{K^*(N \vee J)(n^*)^{-1}\}^{1/4}$ for approximate low rank matrix completion (e.g., Cai and Zhou, 2013; Davenport et al., 2014; Chatterjee et al., 2015), which is expected as the parameter space in these works (nuclear-norm constrained matrices) is larger than that of our setting. Several technical tools are

used to obtain the improved error estimate including a sharp bound on the spectral norm of random matrices which extends a recent result in *Bandeira et al. (2016)* and an upper bound of singular values of Hadamard products of low rank matrices based on a result established in *Horn (1995)*.

Note that the constant κ in Theorem 1 depends on p_{\max}/p_{\min} . Thus, it is most useful when p_{\max}/p_{\min} is bounded by a finite constant that is independent of N and J . In this case, the asymptotic error rate is similar between a uniform sampling and a weighted sampling. In the case where the sampling scheme is far from a uniform sampling, the next theorem provides a finite error bound.

Theorem 2. *Let $\kappa_{2C^2} = \sup_{|x| \leq 2C^2} b''(x)$, $\delta_{C^2} = \frac{1}{2} \inf_{|x| \leq C^2} b''(x)$, $\kappa_{1,b,C,\phi} = 8\delta_{C^2}^{-1}(\phi\kappa_{2C^2})^{1/2} + 16C^2$ and $\kappa_{2,b,C,\phi} = (\phi/C^2) \vee (\phi\kappa_{2C^2})^{1/2}$. Then, there exists a universal constant c such that with probability at least $1 - 2(N + J)^{-1} - (n^*)^{-1}$,*

$$\begin{aligned} & \max_{K^* \leq K \leq K_{\max}} \|\hat{M}^{(K)} - M^*\|_F \\ & \leq p_{\min}^{-1} K_{\max}^{1/2} \{ \kappa_{1,b,C,\phi} (\max_i n_i^*)^{1/2} \vee (\max_j n_j^*)^{1/2} + c(\kappa_{2,b,C,\phi} \log n^* + 2C^2) \log^{1/2}(N + J) \} \end{aligned} \quad (5)$$

for all $N \geq 1, J \geq 1, n^* \geq 6$ and $K_{\max} \geq K^* \geq 1$.

Remark 3. *Theorem 2 provides a finite error bound for the joint maximum likelihood estimator when the number of factors is known to be no greater than K_{\max} . It extends Theorem 1 in several aspects. First, the constants κ_{2C^2} , δ_{C^2} , $\kappa_{1,b,C,\phi}$ and $\kappa_{2,b,C,\phi}$ are made explicit in Theorem 2. In addition, it allows the missing pattern to be far from uniform sampling. To see this, consider the case where $J = N^\alpha$, $p_{\min} = N^{-\beta}$, and $p_{\max}/p_{\min} \leq N^\gamma$ with $\alpha \in (0, 1]$, $\beta \in [0, \alpha)$, $\gamma \in [0, \beta]$ and C is fixed. Roughly, a larger γ suggests a more imbalanced sampling scheme. Then, Theorem 2 implies $(NJ)^{-1/2} \|\hat{M}^{(K^*)} - M^*\|_F = O_p(N^{(\beta+\gamma-\alpha)/2} (K^*)^{1/2})$. Thus, if $\gamma < \alpha - \beta$ and $K^* = o(N^{\alpha-\beta-\gamma})$, the estimator $\hat{M}^{(K^*)}$ is consistent in terms of the scaled Frobenius norm $(NJ)^{-1/2} \|\hat{M}^{(K^*)} - M^*\|_F$.*

Let $u(n, J, K) = v(n, N, J, K) - v(n, N, J, K - 1)$, and let $\sigma_1(M^*) \geq \sigma_2(M^*) \geq \dots \geq \sigma_{K^*+1}(M^*)$ be the non-zero singular values of M^* . Note that due to the inclusion of the intercept term d_j , a non-degenerate M^* is of rank $K^* + 1$. The next theorem provides sufficient conditions on $u(n, N, J, K)$ for consistent model selection.

Theorem 3. *Consider the following asymptotic regime as $N, J \rightarrow \infty$,*

$$C = O(1), p_{\min}^{-1} = O(1), \text{ and } K^* = O(1). \quad (6)$$

If the function u satisfies

$$u(n, N, J, K) = o(\sigma_{K^*+1}^2(M^*)) \text{ and } N \vee J = o(u(n, N, J, K)) \text{ uniformly in } K \text{ as } N, J \rightarrow \infty, \quad (7)$$

then, $\lim_{N, J \rightarrow \infty} \Pr(\hat{K} = K^) = 1$.*

We elaborate on the asymptotic regime (6) and the conditions on $u(n, N, J, K)$ in (7). First, $C = O(1)$ and $K^* = O(1)$ requires that C and the number of factors are bounded as N and J grow. Second, $p_{\min}^{-1} = O(1)$ suggests that the missing pattern is similar to uniform sampling with n^* growing at the order of NJ . Third, $u(n, N, J, K) = o(\sigma_{K^*+1}^2(M^*))$ requires that $u(n, N, J, K)$ is smaller than the gap between non-zero singular values and zero-singular values of M^* . Under this requirement, the probability of underselcting the number of factors is small. Fourth, $N \vee J = o(u(n, J, K))$ requires that $u(n, N, J, K)$ grows in a faster speed than $N \vee J$. This requirement guarantees that with high probability, we do not overselect the number of factors. Fifth, we point out that n is random when there are missing data, and thus $u(n, N, J, K)$ may also be random. In this theorem, we do not allow $u(n, N, J, K)$ to be random as implicitly required by condition (7). A general result allowing a random $u(n, N, J, K)$ is given in Theorem 4 below.

We now discuss concrete choices of the penalty term under the setting of Theorem 3. From (7), informally speaking, $u(n, N, J, K)$ needs to be larger than $N \vee J$ and smaller

than $\sigma_{K^*+1}^2(M^*)$, where $\sigma_{K^*+1}^2(M^*)$ can be viewed as the signal level of the data and $N \vee J$ measures the noise level. As our signal level $\sigma_{K^*+1}^2(M^*)$ is usually unknown in practice, we suggest to choose $u(n, N, J, K)$ to be close to the asymptotic lower bound $N \vee J$, so that it works for a wider range of signal levels. Specifically, consider the case when there is no missing data, i.e., $p_{\min} = 1$. Then a sensible choice would be $u(n, N, J, K) = (N \vee J) \log(N \wedge J)$. This choice is implied by the penalty (3) when there is no missing data. We provide the following corollary of Theorem 3 to establish conditions under which this choice of $u(n, N, J, K)$ leads to consistent model selection. In the presence of missing data, the consistency of the suggested penalty (3) will be established in Corollary 2 as an implication of Theorem 4.

Corollary 1. *Assume that the asymptotic regime (6) holds. Consider $v(n, N, J, K) = K(N \vee J)h(N, J)$ for some function h . If $\lim_{N, J \rightarrow \infty} h(N, J) = \infty$ and $\lim_{N, J \rightarrow \infty} (h(N, J))^{-1}(N \vee J)^{-1}\sigma_{K^*+1}^2(M^*) = \infty$, then $\lim_{N, J \rightarrow \infty} \Pr(\hat{K} = K^*) = 1$. Specifically, suppose that $p_{\min} = 1$. If $(N \vee J) \log(N \wedge J) = o(\sigma_{K^*+1}^2(M^*))$ and we choose $v(n, N, J, K) = K(N \vee J) \log(N \wedge J)$, then $\lim_{N, J \rightarrow \infty} \Pr(\hat{K} = K^*) = 1$.*

The next theorem is a generalisation of Theorem 3 that is established under a more general asymptotic setting.

Theorem 4. *Consider the following asymptotic regime as $N, J \rightarrow \infty$,*

$$C = O(1) \text{ and } (N \wedge J) \log(N + J) = O(n^*/(\log n^*)^2). \quad (8)$$

Also, assume $p_{\min}^{-2} p_{\max} K^(N \vee J) = o(\sigma_{K^*+1}^2(M^*))$. If there exists a (possibly random) sequence $\{\xi_{N, J}\}$ such that $\xi_{N, J} \rightarrow \infty$ in probability as $N, J \rightarrow \infty$, and with probability converg-*

ing to one as $N, J \rightarrow \infty$, the following inequalities hold

$$u(n, N, J, K) \begin{cases} \leq \xi_{N,J}^{-1} p_{\min} \sigma_{K+1}^2(M^*), & \text{if } 1 \leq K \leq K^*, \\ \geq \xi_{N,J}(K^* + 1)(p_{\max}/p_{\min})(N \vee J), & \text{if } K = K^* + 1, \\ \geq \xi_{N,J}(p_{\max}/p_{\min})(N \vee J), & \text{if } K^* + 2 \leq K \leq K_{\max}, \end{cases} \quad (9)$$

where $K_{\max} \geq K^*$ denotes the largest number of factors considered in model selection, where we allow $K_{\max} = \infty$. Then, $\lim_{N,J \rightarrow \infty} \Pr(\hat{K} = K^*) = 1$.

Theorem 4 relaxes the assumptions of Theorem 3 in several aspects. First, it is established under a more general asymptotic regime by allowing K^* to diverge and p_{\min} to decay to zero, as N and J grow. It also allows the missing pattern to be very different from uniform sampling by allowing p_{\max}/p_{\min} to grow. Second, $u(n, N, J, K)$ is allowed to be random as long as (9) holds with high probability. In particular, the model selection consistency of the suggested penalty (3) is established in Corollary 2 below, as an implication of Theorem 4. Third, (9) provides a more specific requirement on $u(n, N, J, K)$. The second and third lines of (9) depend on the true number of factors K^* . In practice, we need to choose $u(n, N, J, K)$ in a way that does not depend on K^* . For example, we may choose $u(n, N, J, K) = (K \wedge K_{\max})(p_{\max}/p_{\min})(N \vee J)h(n, N, J)$ for some sequence $h(n, N, J)$ that tends to infinity in probability as N and J diverge, so that the second and third lines of (9) are satisfied.

Corollary 2. Assume the asymptotic regime (6) holds and $N \vee J = o(\sigma_{K^*+1}^2(M^*))$. Consider $v(n, N, J, K) = K(N \vee J)h(n, N, J)$. If $h(n, N, J) \rightarrow \infty$ in probability as $N, J \rightarrow \infty$ and $(h(n, N, J))^{-1}(N \vee J)^{-1}\sigma_{K^*+1}^2(M^*) \rightarrow \infty$ in probability as $N, J \rightarrow \infty$, then $\lim_{N,J \rightarrow \infty} \Pr(\hat{K} = K^*) = 1$. In particular, if we choose $v(n, N, J, K) = K(N \vee J) \log(n/(N \vee J))$ as suggested in (3) and assume $(N \vee J) \log(N \wedge J) = o(\sigma_{K^*+1}^2(M^*))$, then $\lim_{N,J \rightarrow \infty} \Pr(\hat{K} = K^*) = 1$.

Remark 4. In Theorems 3 and 4, the dispersion parameter ϕ is assumed known. When ϕ is unknown, we may first fit the largest model with K_{\max} factors to obtain an estimate $\hat{\phi}$, and then select the number of factors using JIC with ϕ replaced by $\hat{\phi}$. Similar model

selection consistency results would still hold. We note that the use of the plug-in estimator for dispersion parameter is common in constructing information criteria for linear models and linear factor models (see, e.g., Bai and Ng, 2002).

4 Numerical Experiments

4.1 Simulation

We use a simulation study to evaluate the proposed JIC, focusing on the suggested choice $v(n, N, J, K) = K(N \vee J) \log(n/(N \vee J))$. Twelve settings are considered as listed in Table 1, given by different combinations of models, N , J , and missing data design. Specifically, we generate data from the logistic and Poisson models given in Examples 1 and 2, respectively. Two combinations of N and J are considered, (1) $N = 1000$ and $J = 100$, and (2) $N = 2000$ and $J = 200$. The true number of factors is set to $K^* = 5$. Moreover, we consider three settings for missing data, including (1) no missing data, (2) uniformly missing, with missing probability $p_{ij} = 0.5$ for all i and j , and (3) non-uniformly missing, with missing probability

$$p_{ij} = \exp(f_{i1}^*) / (1 + \exp(f_{i1}^*))$$

that depends on the value of the first factor.

For each setting, 100 independent simulations are run, with the true factor values F_i^* generated from a K^* -variate truncated normal distribution which truncates a K^* -variate standard normal vector in a ball with center at the origin and radius $2\sqrt{2}$. The true factor values thus satisfy $(\|F_i^*\|^2 + 1)^{\frac{1}{2}} \leq 3$. The true parameter vectors $(d_j^*, (A_j^*)^T)^T$ for manifest variables are generated from a (K^*+1) -variate truncated normal distribution which truncates a $(K^* + 1)$ -variate standard normal vector in a ball with center at the origin and radius 3. Under the constraint for F_i , we have the bounds $0.06 \leq p_{ij} \leq 0.94$ for the non-uniform missing setting. For each dataset under each setting, we use the proposed JIC to select K

Setting	Model	N	J	Missing	Correct	Under-selection
1	Logistic	1000	100	No	100	0
2				Uniform	98	2
3				Non-uniform	86	14
4		2000	200	No	100	0
5				Uniform	100	0
6				Non-uniform	100	0
7	Poisson	1000	100	No	100	0
8				Uniform	100	0
9				Non-uniform	100	0
10		2000	200	No	100	0
11				Uniform	100	0
12				Non-uniform	100	0

Table 1: Twelve settings of simulation study and simulation results based on 100 independent replications for each setting.

from the candidate set $\{4, 5, 6\}$ and the constraint constant C in (2) is set to be 4.

Our results are given in Table 1, where for each setting we show the numbers of times among the 100 independent replications that the number of factors is correctly selected and under-selected. Note that over-selection is never observed in these simulations. For the logistic model, we see that under-selection only happens to settings 2 and 3, where both N and J are relatively small and there exist missing data. When N and J are relatively larger or when there is no missing data, the proposed JIC always correctly selects the number of factors. In addition, under-selection is more likely to happen under setting 3 where data are non-uniformly missing. This result is consistent with Theorem 2, where the term $1/p_{min}$ in the error bound (5) for $\hat{M}^{(K)}$ suggests that model estimation tends to be more challenging for smaller p_{min} . Inaccurate estimation further leads to inaccurate model selection. For the Poisson model, the number of factors is always correctly specified under all six settings. These results suggest that the proposed JIC performs well under most of our simulation settings, especially for count data from a Poisson factor model.

4.2 Application to Eysenck’s Personality Questionnaire

We apply the proposed JIC to a dataset based on the revised Eysenck’s personality questionnaire (Eysenck et al., 1985), a personality inventory that has been widely used in clinics and research. This questionnaire is designed to measure three personality traits, including extraversion, neuroticism, and psychoticism. We refer the readers to Eysenck et al. (1985) for the characteristics of these personality traits. We analyse all the items from the questionnaire, except for lie scale items that are used to guard against various concerns about response style. There are 79 items in total, each with “Yes” and “No” response options. An example item is “Do you often need understanding friends to cheer you up?”. Among the 79 items, 32, 23, and 24 items are designed to measure psychoticism, extraversion, and neuroticism, respectively. The determination of the number of latent factors underlying this personality questionnaire is important in personality psychology. Specifically, the number of factors underlying the revised Eysenck’s personality questionnaire has been studied using data from many countries (Barrett et al., 1998).

Here, we analyse the female UK normative sample data (Eysenck et al., 1985) for the questionnaire, for which the sample size is 824 and there is no missing values. The dataset has been analysed in Chen et al. (2019a) using the same model given in Example 1 above. Using a cross-validation approach, Chen et al. (2019a) find three factors. We now explore the dimensionality of the data using the proposed JIC. Specifically, we consider possible choices of $K = 1, 2, 3, 4$, and 5 . Following the previous discussion, the penalty term in the JIC is set to $K(N \vee J) \log(n/(N \vee J))$, where $n = NJ$, $N = 824$, and $J = 79$.

The results are given in Table 2. Specifically, the three-factor model achieves the minimum value of JIC among the five candidate choices of K , suggesting that three factors should be chosen. This result is consistent with the design of the questionnaire, as well as the cross-validation result in Chen et al. (2019a). On the other hand, it is also worth noting that the JIC value of the two-factor model is quite close to that of the three-factor model, suggesting that the signal in the third factor may not be very strong.

K	1	2	3	4	5
Deviance	63263	57683	53883	51225	48812
Penalty	3600	7201	10801	14402	18002
JIC	66864	64884	64684	65627	66814

Table 2: JIC results for an application to the revised Eysenck’s personality questionnaire. The rows named “deviance”, “penalty”, and “JIC” show the values of $-2l_K(\hat{F}_1, \dots, \hat{F}_N, \hat{A}_1, \hat{d}_1, \dots, \hat{A}_J, \hat{d}_J)$, $v(n, N, J, K)$, and JIC, respectively, for models with different values of K .

5 Further Discussion

As shown in Section 3, there is a wide range of penalties for guaranteeing the selection consistency of JIC. Among these choices, $v(n, N, J, K) = K(N \vee J) \log(n/(N \vee J))$ is close to the lower bound. This penalty is suggested when the signal strength of factors is unknown, for being able to detect weak factors. According to simulation study and real data analysis in Section 4, this penalty choice performs well for both the logistic and Poisson factor models. On the other hand, if one is only interested in detecting strong factors, then a larger penalty may be chosen based on prior information about the signal strength of the factors.

When our model (1) takes the form of a Gaussian density and there is no missing data, then the proposed JIC and its theory are consistent with the results of Bai and Ng (2002) for high-dimensional linear factor models. In this sense, the current work substantially extends the work of Bai and Ng (2002) by considering non-linear factor models **and allowing a general setting for missing values**.

Although we focus on generalised latent factor models with an exponential-family link function, the proposed JIC is applicable to other models, for example, a probit factor model for binary data that replaces the logistic link by a probit link in Example 1. The consistency results are likely to hold under similar conditions, for a wider range of models. This extension is left for future investigation.

A Proof of Theoretical Results

A.1 Proof of Theorem 1 and Theorem 2

We will present the proof of Theorem 2 first and then that of Theorem 1 because the former is more general than the latter. The proof of Theorem 2 is based on the following two lemmas, whose proof will be provided later in the supplementary material.

Let $G_i = (1, F_i^T)^T$ and $B_j = (d_j, A_j^T)^T$, then $m_{ij} = G_i^T B_j$. Define $\mathcal{M}_r = \{M = (m_{ij})_{1 \leq i \leq N, 1 \leq j \leq J} : m_{ij} = G_i^T B_j : G_i, B_j \in \mathbb{R}^r, \|G_i\| \leq C, \|B_j\| \leq C \text{ for all } 1 \leq i \leq N, 1 \leq j \leq J\}$, then $M^* \in \mathcal{M}_{K^*+1}$. Let $r^* = K^* + 1$ under Assumption 2. Also, let $l(M, Y, \Omega)$ denote the log-likelihood function where $\Omega = (\omega_{ij})_{1 \leq i \leq N, 1 \leq j \leq J}$.

Lemma 1. *For all $M \in \mathcal{M}_r$,*

$$\begin{aligned} & \phi\{l(M, Y, \Omega) - l(M^*, Y, \Omega)\} \\ & \leq (r + r^*)^{1/2} \left\{ \|Z \circ \Omega\|_2 + 2\delta_{C^2} C^2 \|Q\|_2 \right\} \|M - M^*\|_F - \delta_{C^2} p_{\min} \|M - M^*\|_F^2 \end{aligned} \quad (10)$$

where $Z = (z_{ij})_{1 \leq i \leq N, 1 \leq j \leq J}$, $z_{ij} = y_{ij} - b'(m_{ij}^*)$, and ‘ \circ ’ denotes the matrix Hadamard product.

Lemma 2. *There is a universal constant $c > 0$ such that*

$$\Pr \left(\|\Omega - P\|_2 \geq 4(\max_i n_i^*)^{1/2} \vee (\max_j n_j^*)^{1/2} + c \log^{1/2}(N + J) \right) \leq (N + J)^{-1}. \quad (11)$$

Lemma 3. *Let $V = (v_{ij})_{1 \leq i \leq N, 1 \leq j \leq J}$ be a random matrix with independent and centered entries. In addition, assume v_{ij} s are sub-exponential random variables with parameters $\nu, \alpha > 0$. That is, $E(e^{\lambda v_{ij}}) \leq e^{\lambda^2 \nu^2 / 2}$ for all $|\lambda| < 1/\alpha$. Then, there exists a universal constant $c > 0$ such that with probability at least $1 - (N + J)^{-1} - (n^*)^{-1}$,*

$$\|V \circ \Omega\|_2 \leq 4 \max_{ij} \{E(z_{ij}^2)\}^{1/2} (\max_i n_i^*)^{1/2} \vee (\max_j n_j^*)^{1/2} + c(\alpha \vee \nu) \log n^* \log^{1/2}(N + J) \quad (12)$$

for all $N \geq 1, J \geq 1$, and $n^* \geq 6$. In particular, under Assumptions 1 and 2, $z_{ij} = y_{ij} -$

$b'(m_{ij}^*)$ is sub-exponential with parameters $\nu^2 = \phi \kappa_{2C^2} = \phi \sup_{|x| \leq 2C^2} b''(x)$ and $\alpha = \phi/C^2$, and there is a universal constant $c > 0$ such that with probability at least $1 - (N+J)^{-1} - (n^*)^{-1}$,

$$\|Z \circ \Omega\|_2 \leq 4(\phi \kappa_{2C^2})^{1/2} (\max_i n_{i\cdot}^*)^{1/2} \vee (\max_j n_{\cdot j}^*)^{1/2} + c\{(\phi/C^2) \vee (\phi \kappa_{2C^2})^{1/2}\} \log n^* \log^{1/2}(N+J) \quad (13)$$

for all $N \geq 1, J \geq 1$, and $n^* \geq 6$.

Remark 5. The constant 4 in the first term of the right-hand side of (12) can be improved to $2\sqrt{2} + \epsilon$ for any $\epsilon > 0$ with the constant c replaced by an ϵ -dependent constant c_ϵ . The logarithm term can be improved if Z is further assumed sub-Gaussian or bounded. We keep the current form which is sharp enough for our problem.

Proof of Theorem 2. By the definition of $\hat{M}^{(K)}$ and $K \geq K^*$, we have $\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(M^*, Y, \Omega)\} \geq 0$. Apply Lemma 1 with $M = \hat{M}^{(K)}$, $r = K + 1$ and combine it with $\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(M^*, Y, \Omega)\} \geq 0$. We obtain that for every $K \geq K^*$,

$$\|\hat{M}^{(K)} - M^*\|_F \leq p_{\min}^{-1} (K + K^* + 2)^{1/2} \{\delta_{C^2}^{-1} \|Z \circ \Omega\|_2 + 2C^2 \|Q\|_2\}. \quad (14)$$

Thus,

$$\max_{K^* \leq K \leq K_{\max}} (\|\hat{M}^{(K)} - M^*\|_F) \leq 2p_{\min}^{-1} K_{\max}^{1/2} \{\delta_{C^2}^{-1} \|Z \circ \Omega\|_2 + 2C^2 \|Q\|_2\}, \quad (15)$$

where we used the fact that $K + K^* + 2 \leq 2(K_{\max} + 1) \leq 4K_{\max}$ for $K_{\max} \geq 1$. Apply Lemma 2 and Lemma 3 to obtain an upper bound of the right-hand side of the above inequality and

simplify it. We arrive at

$$\begin{aligned}
& \max_{K^* \leq K \leq K_{\max}} (\|\hat{M}^{(K)} - M^*\|_F) \\
& \leq 2p_{\min}^{-1}(K_{\max})^{1/2} [\{4\delta_{C^2}^{-1}(\phi\kappa_{2C^2})^{1/2} + 8C^2\}(\max_i n_{i.}^*)^{1/2} \vee (\max_j n_{.j}^*)^{1/2} \\
& \quad + c\{(\phi/C^2) \vee (\phi\kappa_{2C^2})^{1/2} \log n^* + 2C^2\} \log^{1/2}(N+J)] \\
& = p_{\min}^{-1}(K_{\max})^{1/2} \{\kappa_{1,b,C,\phi}(\max_i n_{i.}^*)^{1/2} \vee (\max_j n_{.j}^*)^{1/2} + 2c(\kappa_{2,b,C,\phi} \log n^* + 2C^2) \log^{1/2}(N+J)\}
\end{aligned} \tag{16}$$

where we recall that $\kappa_{1,b,C,\phi} = 8\delta_{C^2}^{-1}(\phi\kappa_{2C^2})^{1/2} + 16C^2$ and $\kappa_{2,b,C,\phi} = (\phi/C^2) \vee (\phi\kappa_{2C^2})^{1/2}$. This completes our proof. \square

Proof of Theorem 1. Note that $\max_i n_{i.}^* \leq p_{\max}J$ and $\max_j n_{.j}^* \leq p_{\max}N$. Thus, (5) is simplified to

$$\max_{K^* \leq K \leq K_{\max}} (\|\hat{M}^{(K)} - M^*\|_F) \leq \kappa(K_{\max})^{1/2} \{p_{\min}^{-1/2}(N \vee J)^{1/2} + p_{\min}^{-1} \log n^* \log^{1/2}(N+J)\} \tag{17}$$

for some κ depending on C, b, ϕ and p_{\max}/p_{\min} . Because $p_{\min} = (p_{\min}/p_{\max})p_{\max} \geq (p_{\min}/p_{\max})n^*/(NJ)$, the above inequality implies

$$\begin{aligned}
& \max_{K^* \leq K \leq K_{\max}} (\|\hat{M}^{(K)} - M^*\|_F) \\
& \leq \kappa(K_{\max})^{1/2} \{(n^*/(NJ))^{-1/2}(N \vee J)^{1/2} + (n^*/(NJ))^{-1} \log(n^*) \log^{1/2}(N+J)\}
\end{aligned} \tag{18}$$

with a possibly different κ that also depends on C, b , and ϕ . Multiplying both sides by $(NJ)^{-1/2}$ and simplifying it, we arrive at

$$\begin{aligned}
& \max_{K^* \leq K \leq K_{\max}} \{(NJ)^{-1/2} \|\hat{M}^{(K)} - M^*\|_F\} \\
& \leq \kappa K_{\max}^{1/2} \left[\{(N \vee J)/n^*\}^{1/2} + \{(NJ)^{1/2} \log^{1/2}(N+J)\} (n^*)^{-1} \log n^* \right].
\end{aligned} \tag{19}$$

Note that for $n^*/(\log n^*)^2 \geq (N \wedge J) \log(N+J)$, $\{(N \vee J)/n^*\}^{1/2} \geq \{(NJ)^{1/2} \log^{1/2}(N+J)\} (n^*)^{-1} \log n^*$.

$J)\}(n^*)^{-1} \log n^*$, and the above inequality is simplified as

$$\max_{K^* \leq K \leq K_{\max}} \{(NJ)^{-1/2} \|\hat{M}^{(K)} - M^*\|_F\} \leq 2\kappa \{K_{\max}(N \vee J)/n^*\}^{1/2}. \quad (20)$$

This completes the proof. \square

A.2 Proof of Theorem 3, Theorem 4, and Corollary 2

The proofs of Theorem 3 and Theorem 4 are based on the following three supporting lemmas, whose proofs are given in the supplementary material. We start by recalling $u(n, N, J, K) = v(n, N, J, K) - v(n, N, J, K - 1)$ and defining $R = 4(p_{\min} \delta_{C^2})^{-1} \{\|Z \circ \Omega\|_2 + 2\delta_{C^2} C^2 \|Q\|_2\}^2$.

Lemma 4. *If $u(\cdot)$ satisfies*

$$\lim_{N, J \rightarrow \infty} \Pr \left(u(n, N, J, K^* + 1) > 2\phi^{-1}(K^* + 1)R \right) = 1 \quad (21)$$

and

$$\lim_{N, J \rightarrow \infty} \Pr \left(\inf_{K^* + 2 \leq K \leq K_{\max}} u(n, N, J, K) > 2\phi^{-1}R \right) = 1, \quad (22)$$

then

$$\lim_{N, J \rightarrow \infty} \Pr(\hat{K} > K^*) = 0, \quad (23)$$

for $K_{\max} \geq K^* \geq 1$.

Lemma 5. *If*

$$\lim_{N, J \rightarrow \infty} \Pr \left(4(\delta_{C^2} p_{\min})^{-1} K^* R \leq \sigma_{K^*+1}^2(M^*) \right) = 1, \quad (24)$$

and $u(\cdot)$ satisfies

$$\lim_{N, J \rightarrow \infty} \Pr \left(u(n, N, J, K) < \phi^{-1} \delta_{C^2} p_{\min} \sigma_{K+1}^2(M^*) \text{ for all } 1 \leq K \leq K^* \right) = 1 \quad (25)$$

then $\lim_{N, J \rightarrow \infty} \Pr(\hat{K} < K^*) = 0$ for $K^* \geq 1$.

Lemma 6. *Under the asymptotic regime (8), $R = O_p(p_{\max}/p_{\min}(N \vee J))$.*

In the rest of the section, we provide the proof of Theorem 4 first and then the proof of Theorem 3 because the former is more general than the latter.

Proof of Theorem 4. In the proof, we will verify that conditions of Theorem 4 ensure conditions in Lemma 4 and Lemma 5. We start with verifying conditions in Lemma 4. According to the second line of (9),

$$\begin{aligned}
& \lim_{N, J \rightarrow \infty} \Pr \left(u(n, N, J, K^* + 1) > 2\phi^{-1}(K^* + 1)R \right) \\
& \geq \liminf_{N, J \rightarrow \infty} \Pr \left(\xi_{N, J}(K^* + 1)(p_{\max}/p_{\min})(N \vee J) > 2\phi^{-1}(K^* + 1)R \right) \\
& \geq \liminf_{N, J \rightarrow \infty} \Pr \left(\xi_{N, J}(p_{\max}/p_{\min})(N \vee J) > 2\phi^{-1}R \right) \\
& = 1,
\end{aligned} \tag{26}$$

where the last line is obtained according to Lemma 6 and that $\xi_{N, J} \rightarrow \infty$ in probability.

Similarly,

$$\begin{aligned}
& \lim_{N, J \rightarrow \infty} \Pr \left(\inf_{K^* + 2 \leq K \leq K_{\max}} u(n, N, J, K) > 2\phi^{-1}R \right) \\
& \geq \liminf_{N, J \rightarrow \infty} \Pr \left(\xi_{N, J}(p_{\max}/p_{\min})(N \vee J) > 2\phi^{-1}R \right) \\
& = 1.
\end{aligned} \tag{27}$$

Thus, conditions of Lemma 4 are verified and we obtain

$$\lim_{N, J \rightarrow \infty} \Pr(\hat{K} > K^*) = 0. \tag{28}$$

Next, we verify conditions of Lemma 5. According to Lemma 6 and the assumption $p_{\min}^{-2}p_{\max}K^*(N \vee J) = o(\sigma_{K^*+1}^2(M^*))$, we have

$$4(\delta_{C^2}p_{\min})^{-1}K^*R = O_p(p_{\min}^{-2}p_{\max}K^*(N \vee J)) = o_p(\sigma_{K^*}^2(M^*)). \tag{29}$$

Thus,

$$\lim_{N,J \rightarrow \infty} \Pr \left(4(\delta_{C^2} p_{\min})^{-1} K^* R \leq \sigma_{K^*}^2(M^*) \right) = 1. \quad (30)$$

In addition, according to the first line of (9),

$$\begin{aligned} & \lim_{N,J \rightarrow \infty} \Pr \left(u(n, N, J, K) < \phi^{-1} \delta_{C^2} p_{\min} \sigma_{K+1}^2(M^*) \text{ for } 1 \leq K \leq K^* \right) \\ & \geq \liminf_{N,J \rightarrow \infty} \Pr \left(\xi_{N,J}^{-1} p_{\min} \sigma_{K^*+1}^2(M^*) < \phi^{-1} \delta_{C^2} p_{\min} \sigma_{K+1}^2(M^*) \text{ for all } K \right) \\ & \geq \liminf_{N,J \rightarrow \infty} \Pr \left(\xi_{N,J}^{-1} < \phi^{-1} \delta_{C^2} \right) \\ & = 1. \end{aligned} \quad (31)$$

From (30) and (31), conditions of Lemma 5 are verified and thus

$$\lim_{N,J \rightarrow \infty} \Pr(\hat{K} < K^*) = 0. \quad (32)$$

We complete the proof by combining (28) and (32). \square

Proof of Theorem 3. First note that the existence of u satisfying (7) implies $N \vee J = o(\sigma_{K^*+1}^2(M^*))$, which further implies $p_{\min}^{-2} p_{\max} K^* (N \vee J) = o(\sigma_{K^*+1}^2(M^*))$ under the asymptotic regime $p_{\min}^{-1} = O(1)$, $K^* = O(1)$. Thus, the assumption about the singular value of M^* in Theorem 4 is verified. Also, $p_{\min}^{-1} = O(1)$ implies that $(N \wedge J) \log(N+J) = o(n^*/(\log n^*)^2)$. Thus, (8) is verified.

We proceed to verify that u satisfies (9) in Theorem 4. We note that $p_{\min}^{-1} = O(1)$, $K^* = O(1)$ and u satisfies (7) implies that there exists $\xi_{N,J} \rightarrow \infty$ satisfying

$$u(n, N, J, K) \leq \xi_{N,J}^{-1} p_{\min} \sigma_{K^*+1}^2(M^*) \text{ for all } K, \quad (33)$$

$$u(n, N, J, K) \geq \xi_{N,J} (p_{\max}/p_{\min}) (N \vee J) \text{ for all } K, \quad (34)$$

and

$$u(n, N, J, K^* + 1) \geq \xi_{N,J}(K^* + 1)(p_{\max}/p_{\min})(N \vee J). \quad (35)$$

Note that $\sigma_{K+1}^2(M^*) \geq \sigma_{K^*+1}^2(M^*)$ for $K \leq K^*$. Thus, (33) implies the first line of (9); (35) implies the second line of (9); and (34) implies the last line of (9). This verifies (9) and completes the proof. \square

Proof of Corollary 2. Under the asymptotic regime (6) and $N \vee J = o(\sigma_{K^*+1}^2(M^*))$, (8) and $p_{\min}^{-2} p_{\max} K^* (N \vee J) = o(\sigma_{K^*+1}^2(M^*))$ are already verified in the proof of Theorem 3. We proceed to verify (9).

From the conditions on $h(n, N, J)$, there exists a sequence $\xi_{N,J}$ (possibly depending on $h(n, N, J)$) such that $\xi_{N,J} \rightarrow \infty$ in probability and

$$\xi_{N,J} < h(n, N, J)(p_{\min}/p_{\max})(K^* + 1)^{-1} \text{ and } \xi_{N,J} \leq (h(n, N, J))^{-1}(N \vee J)^{-1} p_{\min} \sigma_{K^*+1}^2(M^*). \quad (36)$$

Also, note that $u(n, N, J, K) = v(n, N, J, K) - v(n, N, J, K - 1) = (N \vee J)h(n, N, J)$. It is not hard to verify (36) implies (9), and, thus Theorem 4 applies.

We proceed to the proof of the ‘in particular’ part. Note that by definition $E(n) = n^*$ and $\text{Var}(n) = \sum_i \sum_j p_{ij}(1 - p_{ij}) \leq \sum_i \sum_j p_{ij} = n^*$, which implies $\lim_{N,J \rightarrow \infty} \Pr(n > 2n^* \text{ or } n < n^*/2) = 0$ and further implies

$$\lim_{N,J \rightarrow \infty} \Pr \left(n/(N \vee J) \geq 2n^*/(N \vee J) \text{ or } n/(N \vee J) \leq n^*/\{2(N \vee J)\} \right) = 0.$$

Note that in this part, $h(n, N, J) = \log(n/(N \vee J))$. Also, $\log(n^*/\{2(N \vee J)\}) \rightarrow \infty$. Thus, $h(n, N, J) \rightarrow \infty$ in probability. In addition, on the event $n/(N \vee J) \leq 2n^*/(N \vee J)$, $(h(n, N, J))^{-1}(N \vee J)^{-1} \sigma_{K^*+1}^2(M^*) \geq \log(2n^*/(N \vee J))(N \vee J)^{-1} \sigma_{K^*+1}^2(M^*)$. The right-hand-side of this inequality tend to infinity under the assumptions of the Corollary. This implies $(h(n, N, J))^{-1}(N \vee J)^{-1} \sigma_{K^*+1}^2(M^*) \rightarrow \infty$ in probability. \square

A.3 Proof of Supporting Lemmas

Proof of Lemma 1. By definition,

$$\begin{aligned}
& \phi\{l(M, Y, \Omega) - l(M^*, Y, \Omega)\} \\
&= \sum_{ij} \omega_{ij} \{y_{ij}m_{ij} - b(m_{ij}) - y_{ij}m_{ij}^* + b(m_{ij}^*)\} \\
&= \sum_{i,j} (y_{ij} - b'(m_{ij}^*))(m_{ij} - m_{ij}^*)\omega_{ij} - \sum_{ij} \{b(m_{ij}) - b(m_{ij}^*) - b'(m_{ij}^*)(m_{ij} - m_{ij}^*)\}\omega_{ij}.
\end{aligned} \tag{37}$$

In the rest of the proof, we derive upper bounds for each term on the right-hand-side of the above display. For the first term $\sum_{i,j} (y_{ij} - b'(m_{ij}^*))(m_{ij} - m_{ij}^*)\omega_{ij}$, we write it as

$$\sum_{i,j} (y_{ij} - b'(m_{ij}^*))(m_{ij} - m_{ij}^*)\omega_{ij} = \langle Z \circ \Omega, M - M^* \rangle, \tag{38}$$

where $\langle A, B \rangle = \text{tr}(A^T B)$ denotes the matrix inner product. Recall the following inequality in linear algebra: $|\langle A, B \rangle| \leq \|A\|_2 \|B\|_* \leq \sqrt{\text{rank}(B)} \|A\|_2 \|B\|_F$ for any two matrices A and B . Applying this fact to the above display, we obtain

$$\left| \sum_{i,j} (y_{ij} - b'(m_{ij}^*))(m_{ij} - m_{ij}^*)\omega_{ij} \right| \leq \{\text{rank}(M - M^*)\}^{1/2} \|Z \circ \Omega\|_2 \|M - M^*\|_F. \tag{39}$$

Notice that $\text{rank}(M - M^*) \leq r + r^*$ for $M \in \mathcal{M}_r$. Thus, the above inequality implies

$$\left| \sum_{i,j} (y_{ij} - b'(m_{ij}^*))(m_{ij} - m_{ij}^*)\omega_{ij} \right| \leq (r + r^*)^{1/2} \|Z \circ \Omega\|_2 \|M - M^*\|_F. \tag{40}$$

We proceed to the analysis of the second term $\sum_{ij} \{b(m_{ij}) - b(m_{ij}^*) - b'(m_{ij}^*)(m_{ij} - m_{ij}^*)\}\omega_{ij}$. Note that for $M \in \mathcal{M}_r$, $|m_{ij}| \leq \|B_i\| \|G_j\| \leq C^2$. Similarly, $|m_{ij}^*| \leq C^2$. Thus, for any $\tilde{m}_{ij} = tm_{ij}^* + (1-t)m_{ij}$ and $t \in (0, 1)$, $|\tilde{m}_{ij}| \leq C^2$. Recall the definition of $\delta_{C^2} = \inf_{|x| \leq C^2} b''(x)$.

Then, $\frac{1}{2}b''(\tilde{m}_{ij}) \geq \delta_{C^2}$. This implies

$$\begin{aligned}
& \sum_{ij} \{b(m_{ij}) - b(m_{ij}^*) - b'(m_{ij}^*)(m_{ij} - m_{ij}^*)\} \omega_{ij} \\
&= \sum_{ij} \frac{1}{2} b''(\tilde{m}_{ij}) (m_{ij} - m_{ij}^*)^2 \omega_{ij} \\
&\geq \delta_{C^2} \sum_{ij} (m_{ij} - m_{ij}^*)^2 \omega_{ij}.
\end{aligned} \tag{41}$$

Note that

$$\begin{aligned}
& \sum_{ij} (m_{ij} - m_{ij}^*)^2 \omega_{ij} \\
&= \sum_{ij} (m_{ij} - m_{ij}^*)^2 (\omega_{ij} - p_{ij}) + \sum_{ij} p_{ij} (m_{ij} - m_{ij}^*)^2 \\
&\geq \langle (M - M^*) \circ (M - M^*), Q \rangle + p_{\min} \|M - M^*\|_F^2 \\
&\geq -\|(M - M^*) \circ (M - M^*)\|_* \|Q\|_2 + p_{\min} \|M - M^*\|_F^2.
\end{aligned} \tag{42}$$

where we define $Q = \Omega - P$ and $P = (p_{ij})_{1 \leq i \leq N, 1 \leq j \leq J}$. The next lemma is helpful for bounding matrix norms involving Hadamard products, whose proof is given later this section.

Lemma 7. *For $M \in \mathcal{M}_r$, $\|(M - M^*) \circ (M - M^*)\|_* \leq 2C^2(r + r^*)^{1/2} \|M - M^*\|_F$.*

Remark 6. *The proof of Lemma 7 utilizes the property that $m_{ij} = B_i^T G_j$ with $\|B_i\|, \|G_j\| \leq C$ and combine it with a result in Horn (1995). This improves the estimate in Chen et al. (2019b) where $|m_{ij}| \leq C^2$ is directly used to derive an upper bound $2C^2(r + r^*) \|M - M^*\|_F$. Comparing with this bound, the above lemma provide a sharper bound in the order of $r + r^*$.*

Applying Lemma 7 to (42) and combine it with (41), we obtain

$$\begin{aligned}
& \sum_{ij} \{b(m_{ij}) - b(m_{ij}^*) - b'(m_{ij}^*)(m_{ij} - m_{ij}^*)\} \omega_{ij} \\
&\geq \delta_{C^2} \{p_{\min} \|M - M^*\|_F^2 - 2C^2(r + r^*)^{1/2} \|M - M^*\|_F \|Q\|_2\}.
\end{aligned} \tag{43}$$

We complete the proof by combining the above display with (37) and (40). \square

Proof of Lemma 7. Let $\tilde{B}_i = (B_i^T, -(B_i^*)^T)^T$ and $\tilde{G}_j = (G_j^T, -(G_j^*)^T)^T$. Then, $\tilde{B}_i, \tilde{G}_j \in \mathbb{R}^{r+r^*}$, $\|\tilde{B}_i\|, \|\tilde{G}_j\| \leq \sqrt{2}C$, and $m_{ij} - m_{ij}^* = \tilde{B}_i^T \tilde{G}_j$ for all i, j .

On the other hand, Theorem 2 in Horn (1995) states that, for any $m \times n$ matrices $A = (a_{ij}), B = (b_{ij})$, if $a_{ij} = g_j^T f_i$ for vectors g_j and f_i s. Then,

$$\sum_{i=1}^k \sigma_i(A \circ B) \leq \sum_{i=1}^k \|f_{[i]}\| \|g_{[i]}\| \sigma_i(B) \text{ for } k = 1, \dots, m \wedge n, \quad (44)$$

where $\sigma_i(\cdot)$ denotes the i th largest singular value of a matrix, $\|f_{[1]}\| \geq \|f_{[2]}\| \geq \dots \geq \|f_{[m]}\|$ and $\|g_{[1]}\| \geq \dots \geq \|g_{[n]}\|$ denote the order statistics of $\{\|f_i\|\}_{i=1}^m$ and $\{\|g_j\|\}_{j=1}^n$. Now, we let $k = N \wedge J$, $A = M - M^*$, $B = A$, $f_i = \tilde{B}_i$, $g_j = \tilde{G}_j$ in the above result and note that $\|f_{[i]}\|, \|g_{[j]}\| \leq \sqrt{2}C$ in this case, we obtain

$$\sum_{i=1}^{N \wedge J} \sigma_i((M - M^*) \circ (M - M^*)) \leq \sum_{i=1}^{N \wedge J} 2C^2 \sigma_i(M - M^*) = 2C^2 \|M - M^*\|_*. \quad (45)$$

Noting the left-hand side of the above display equals $\|(M - M^*) \circ (M - M^*)\|_*$. Thus,

$$\|(M - M^*) \circ (M - M^*)\|_* \leq 2C^2 \|M - M^*\|_* \leq 2C^2 (r + r^*)^{1/2} \|M - M^*\|_F. \quad (46)$$

\square

The proofs of Lemmas 2 and 3 are based on the next lemma that provides an upper tail bound for the spectral norm of a large class of random matrices. Its proof mainly combines standard symmetrization and truncation arguments with a recent result by Bandeira et al. (2016) on the spectral norm of symmetric random matrices with independent, centered and symmetric entries.

Lemma 8. *Let $X = (x_{ij})_{1 \leq i \leq N, 1 \leq j \leq J}$ be an $N \times J$ matrix with $E(x_{ij}) = 0$ and $E(x_{ij}^2) < \infty$.*

Then, there is a universal constant $c > 0$ such that for all $t, \lambda \geq 0$

$$\Pr \left(\|X\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t \right) \leq (N + J)e^{-t^2/(c\lambda^2)} + \sum_{i=1}^N \sum_{j=1}^J \Pr(|x_{ij} - x'_{ij}| > \lambda), \quad (47)$$

where we define $\sigma_1 = \max_{1 \leq i \leq N} \{\sum_{j=1}^J E(x_{ij}^2)\}^{1/2}$, $\sigma_2 = \max_{1 \leq j \leq J} \{\sum_{i=1}^N E(x_{ij}^2)\}^{1/2}$, and x'_{ij} is an independent copy of x_{ij} .

Proof of Lemma 8. Let $X' = (x'_{ij})$ which is an independent copy of X and let $\tilde{X} = (\tilde{x}_{ij}) = X - X'$. Then, \tilde{x}_{ij} s have symmetric distribution and are independent. Let $Z = (z_{ij}) = \begin{pmatrix} 0 & \tilde{X} \\ \tilde{X}^T & 0 \end{pmatrix}$. Z is a symmetric $(N+J) \times (N+J)$ random matrix whose entries are independent and symmetric random variables. Define a random matrix $Z(\lambda)$ as the truncated Z ,

$$Z(\lambda) = (z_{ij}(\lambda))_{1 \leq i \leq N, 1 \leq j \leq J} = (z_{ij}I(|z_{ij}| \leq \lambda))_{1 \leq i \leq N, 1 \leq j \leq J}. \quad (48)$$

Then, entries of $Z(\lambda)$ are independent, symmetric random variables and are bounded by λ . Apply Corollary 3.12 in Bandeira et al. (2016) to $Z(\lambda)$, then there exists a universal constant $c > 0$ such that

$$\Pr \left(\|Z(\lambda)\|_2 \geq 2^{3/2} \max_{1 \leq i \leq (N+J)} \left[\sum_{j=1}^{N+J} E\{z_{ij}^2(\lambda)\} \right]^{1/2} + t \right) \leq (N + J)e^{-t^2/(c\lambda^2)} \quad (49)$$

Note that

$$\begin{aligned} \max_{1 \leq i \leq (N+J)} \left[\sum_{j=1}^{N+J} E\{z_{ij}^2(\lambda)\} \right]^{1/2} &\leq \max_{1 \leq i \leq (N+J)} \left\{ \sum_{j=1}^{N+J} E(z_{ij}^2) \right\}^{1/2} \\ &= \max \left[\max_{1 \leq i \leq N} \left\{ \sum_{j=1}^J E(\tilde{x}_{ij}^2) \right\}^{1/2}, \max_{1 \leq j \leq J} \left\{ \sum_{i=1}^N E(\tilde{x}_{ij}^2) \right\}^{1/2} \right] \\ &= \sqrt{2}(\sigma_1 \vee \sigma_2). \end{aligned} \quad (50)$$

Thus,

$$\Pr \left(\|Z(\lambda)\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t \right) \leq (N + J)e^{-t^2/(c\lambda^2)}. \quad (51)$$

On the other hand,

$$\Pr \left(\|Z\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t \right) \leq \Pr \left(\|Z(\lambda)\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t \right) + \Pr \left(\max_{1 \leq i, j \leq N+J} |z_{ij}| > \lambda \right). \quad (52)$$

The above two inequalities together imply

$$\Pr \left(\|Z\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t \right) \leq (N + J)e^{-t^2/(c\lambda^2)} + \Pr \left(\max_{1 \leq i, j \leq N+J} |z_{ij}| > \lambda \right). \quad (53)$$

Note that $\|Z\|_2 = \|\tilde{X}\|_2$ and $\max_{1 \leq i, j \leq N+J} |z_{ij}| = \max_{1 \leq i \leq N, 1 \leq j \leq J} |\tilde{x}_{ij}|$. From the above inequality, we obtain

$$\Pr \left(\|\tilde{X}\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t \right) \leq (N + J)e^{-t^2/(c\lambda^2)} + \Pr \left(\max_{1 \leq i \leq N, 1 \leq j \leq J} |\tilde{x}_{ij}| > \lambda \right). \quad (54)$$

With a union bound, we further get

$$\Pr(\|\tilde{X}\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t) \leq (N + J)e^{-t^2/(c\lambda^2)} + \sum_{1 \leq i \leq N, 1 \leq j \leq J} \Pr(|\tilde{x}_{ij}| > \lambda). \quad (55)$$

Recall $\tilde{X} = X - X'$ and the function $I(\|X - X'\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t)$ is convex in X' . Thus, by Jensen's inequality,

$$\Pr(\|X\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t) \leq \Pr(\|X - X'\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t) = \Pr(\|\tilde{X}\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t). \quad (56)$$

This, together with (55) completes the proof. \square

Proof of Lemma 2. Let ω'_{ij} be an independent copy of ω_{ij} , then $|\omega'_{ij} - p_{ij} - (\omega_{ij} - p_{ij})| \leq 1$. In addition, $E(\omega_{ij} - p_{ij})^2 = p_{ij}(1 - p_{ij}) \leq p_{ij}$. Thus, $\max_i \{\sum_j E(\omega_{ij} - p_{ij})^2\}^{1/2} \leq$

$\max_i (\sum_j p_{ij})^{1/2} = (\max_i n_i^*)^{1/2}$ and $\max_j \{\sum_i E(\omega_{ij} - p_{ij})^2\}^{1/2} \leq (\max_j n_j^*)^{1/2}$.

Choose $\lambda = 1$ and apply Lemma 8 to $\Omega - P$, we obtain that for all $t \geq 0$,

$$\Pr(\|\Omega - P\|_2 \geq 4(\max_i n_i^* \vee \max_j n_j^*)^{1/2} + t) \leq (N + J)e^{-t^2/c}. \quad (57)$$

Let $t = (2c \log(N + J))^{1/2}$ in the above inequality, we obtain

$$\Pr(\|\Omega - P\|_2 \geq 4(\max_i n_i^* \vee \max_j n_j^*)^{1/2} + (2c \log(N + J))^{1/2}) \leq (N + J)^{-1}. \quad (58)$$

We complete the proof by noting that $(2c)^{1/2}$ is still a universal constant. \square

Proof of Lemma 3. Apply Lemma 8 to $V \circ \Omega$, we obtain that for all $t, \lambda \geq 0$,

$$\Pr(\|V \circ \Omega\|_2 \geq 4(\sigma_1 \vee \sigma_2) + t) \leq (N + J)e^{-t^2/(c\lambda^2)} + \sum_{ij} \Pr(|v_{ij}\omega_{ij} - v'_{ij}\omega'_{ij}| \geq \lambda), \quad (59)$$

where (v'_{ij}, ω'_{ij}) is an independent copy of (v_{ij}, ω_{ij}) , $\sigma_1 = \max_i \{\sum_j E(v_{ij}^2 \omega_{ij}^2)\}^{1/2}$ and $\sigma_2 = \max_j \{\sum_i E(v_{ij}^2 \omega_{ij}^2)\}^{1/2}$. We proceed to a detailed analysis of σ_1, σ_2 and the probability $\Pr(|v_{ij}\omega_{ij} - v'_{ij}\omega'_{ij}| \geq \lambda)$. First, a direct calculation gives

$$\sigma_1 = \max_i \left\{ \sum_j p_{ij} E(v_{ij}^2) \right\}^{1/2} \leq (\max_i n_i^*)^{1/2} \max_{ij} \{E(v_{ij}^2)\}^{1/2}. \quad (60)$$

Similarly, $\sigma_2 \leq (\max_j n_j^*)^{1/2} \max_{ij} \{E(v_{ij}^2)\}^{1/2}$. Now we find an upper bound of $\Pr(|v_{ij}\omega_{ij} - v'_{ij}\omega'_{ij}| \geq \lambda)$. Note that

$$\begin{aligned} & \Pr(|v_{ij}\omega_{ij} - v'_{ij}\omega'_{ij}| \geq \lambda) \\ &= p_{ij}^2 \Pr(|v_{ij} - v'_{ij}| \geq \lambda) + 2p_{ij}(1 - p_{ij}) \Pr(|v_{ij}| \geq \lambda) \\ &\leq 3p_{ij} \Pr(|v_{ij} - v'_{ij}| \geq \lambda) \vee \Pr(|v_{ij}| \geq \lambda). \end{aligned} \quad (61)$$

For $\Pr(|v_{ij}| \geq \lambda)$, we use a tail bound for sub-exponential variables

$$\Pr(|v_{ij}| \geq \lambda) \leq 2e^{-\lambda^2/(2\nu^2)} \vee e^{-\lambda/(2\alpha)}. \quad (62)$$

Similarly, noting that $v_{ij} - v'_{ij}$ is also sub-exponential with parameters $2\nu^2, \alpha$, we have

$$\Pr(|v_{ij} - v'_{ij}| \geq \lambda) \leq 2e^{-\lambda^2/(4\nu^2)} \vee e^{-\lambda/(2\alpha)}. \quad (63)$$

Combining the above two inequalities with (61), we have

$$\Pr(|v_{ij}\omega_{ij} - v'_{ij}\omega'_{ij}| \geq \lambda) \leq 6p_{ij}e^{-\lambda^2/(4\nu^2)} \vee e^{-\lambda/(2\alpha)}. \quad (64)$$

Combining the above inequality with (59), we arrive at

$$\begin{aligned} \Pr(\|V \circ \Omega\|_2 \geq 4 \max_{ij} \{E(v_{ij}^2)\}^{1/2} (\max_i n_i^*)^{1/2} \vee (\max_j n_j^*)^{1/2} + t) \\ \leq (N+J)e^{-t^2/(c\lambda^2)} + 6e^{-\lambda^2/(4\nu^2)} \vee e^{-\lambda/(2\alpha)} n^*. \end{aligned} \quad (65)$$

Let $\lambda = 4(\alpha \vee \nu) \log n^*$. It is not hard to verify that $6e^{-\lambda^2/(4\nu^2)} \vee e^{-\lambda/(2\alpha)} n^* \leq (n^*)^{-1}$ for $n^* \geq 6$. Let $t = \lambda\{2c \log(N+J)\}^{1/2}$, we obtain $(N+J)e^{-t^2/(c\lambda^2)} \leq (N+J)^{-1}$. Combining the above inequalities with (65), we obtain that with probability at least $1 - (N+J)^{-1} - (n^*)^{-1}$,

$$\begin{aligned} \|V \circ \Omega\|_2 \\ \leq 4 \max_{ij} \{E(v_{ij}^2)\}^{1/2} (\max_i n_i^*)^{1/2} \vee (\max_j n_j^*)^{1/2} + 4\sqrt{2}c^{1/2}(\alpha \vee \nu) \log n^* \log^{1/2}(N+J). \end{aligned} \quad (66)$$

This completes the proof of inequality (12) (note that $4\sqrt{2}c^{1/2}$ is also a universal constant).

We proceed to prove the ‘in particular’ part of the lemma. For each $z_{ij} = y_{ij} - b'(m_{ij}^*)$, its second moment is $E(z_{ij}^2) = \phi b''(m_{ij}^*) \leq \phi \kappa_2 C^2$. In addition, its moment generating function is $E(e^{\lambda z_{ij}}) = \exp[\phi^{-1}\{b(m_{ij}^* + \lambda\phi) - b(m_{ij}^*)\} - \lambda b'(m_{ij}^*)] = \exp\{\phi b''(m_{ij}^* + \tilde{\lambda}\phi)\lambda^2/2\}$ for some $|\tilde{\lambda}| \leq |\lambda|$. Since $|m_{ij}^*| \leq C^2$ by assumption, we can see that for $|\lambda| \leq C^2/\phi$, $|m_{ij}^* + \tilde{\lambda}\phi| \leq 2C^2$

and thus $E(e^{\lambda z_{ij}}) \leq \exp\{\kappa_{2C^2}\phi\lambda^2/2\}$ for all $|\lambda| \leq C^2/\phi$. This implies that z_{ij} is sub-exponential with the parameters $\nu^2 = \phi\kappa_{2C^2}$ and $\alpha = \phi/C^2$. We complete the proof by applying (12) with the above parameters for Z . \square

Proof of Lemma 4. For each $K^* + 1 \leq K \leq K_{\max}$, we first derive an upper bound for $\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(\hat{M}^{(K^*)}, Y, \Omega)\} - (v(n, N, J, K) - v(n, N, J, K^*))$. According to Lemma 1,

$$\begin{aligned} & \phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(M^*, Y, \Omega)\} \\ & \leq (K + K^* + 2)^{1/2} (\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2) \|\hat{M}^{(K)} - M^*\|_F \\ & \leq 2K^{1/2} (\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2) \|\hat{M}^{(K)} - M^*\|_F. \end{aligned} \tag{67}$$

Combining this with (14) gives

$$\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(M^*, Y, \Omega)\} \leq 4p_{\min}^{-1}K \left\{ \|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2 \right\}^2 = KR. \tag{68}$$

Thus, the penalized log-likelihood satisfies

$$\begin{aligned} & \max_{K^*+1 \leq K \leq K_{\max}} \left[-2l(\hat{M}^{(K)}, Y, \Omega) + v(n, N, J, K) - \{-2l(M^*, Y, \Omega) + v(n, N, J, K^*)\} \right] \\ & \geq \max_{K^*+1 \leq K \leq K_{\max}} \left[-2\phi^{-1}KR + \sum_{l=K^*+1}^K u(n, N, J, l) \right] \end{aligned} \tag{69}$$

It is easy to see that, if the events $u(n, N, J, K^* + 1) > 2\phi^{-1}(K^* + 1)R$ and $u(n, N, J, l) > 2\phi^{-1}R$ happen at the same time for all $K^* + 2 \leq l \leq K_{\max}$, then the right-hand side of the

above inequality is strictly greater than zero. Thus,

$$\begin{aligned}
& \Pr\left(\hat{K} \leq K^*\right) \\
& \geq \Pr\left(\max_{K^*+1 \leq K \leq K_{\max}} \left[-2l(\hat{M}^{(K)}, Y, \Omega) + v(n, N, J, K) - \{-2l(M^*, Y, \Omega) + v(n, N, J, K^*)\}\right] > 0\right) \\
& \geq \Pr\left(u(n, N, J, K^* + 1) > 2\phi^{-1}(K^* + 1)R, \text{ and } \inf_{K^*+2 \leq l \leq K_{\max}} u(n, N, J, K) > 2\phi^{-1}R\right).
\end{aligned} \tag{70}$$

We complete the proof by noting the right-hand side of the above inequality tend to one under the assumptions of the lemma. \square

The proof of Lemma 5 requires the next lemma.

Lemma 9. *If $4(\delta_{C^2} p_{\min})^{-1} K^* R \leq \sigma_{K^*+1}^2(M^*)$, then*

$$\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(\hat{M}^{(K^*)}, Y, \Omega)\} \leq -\frac{1}{2}\delta_{C^2} p_{\min} \left\{ \sum_{l=K+2}^{K^*+1} \sigma_l^2(M^*) \right\} \tag{71}$$

for $0 \leq K \leq K^* - 1$.

Proof of Lemma 9. First, according to Lemma 1, $\hat{M}^{(K)} \in \mathcal{M}_{K+1}$ and $K + 1 \leq K^*$, we have

$$\begin{aligned}
& \phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(M^*, Y, \Omega)\} \\
& \leq (K + K^* + 2)^{1/2} \left\{ \|Z \circ \Omega\|_2 + 2\delta_{C^2} C^2 \|Q\|_2 \right\} \|\hat{M}^{(K)} - M^*\|_F - \delta_{C^2} p_{\min} \|\hat{M}^{(K)} - M^*\|_F^2 \\
& \leq \sup_{M \in \mathcal{M}_{K+1}} \left[2(K^*)^{1/2} \left\{ \|Z \circ \Omega\|_2 + 2\delta_{C^2} C^2 \|Q\|_2 \right\} \|M - M^*\|_F - \delta_{C^2} p_{\min} \|M - M^*\|_F^2 \right].
\end{aligned} \tag{72}$$

Note that the expression inside ‘sup’ is a quadratic function in $\|M - M^*\|_F$. Let $d(M^*, \mathcal{M}_{K+1}) = \inf_{M \in \mathcal{M}_{K+1}} \|M - M^*\|_F$. From properties of a quadratic function, if $d(M^*, \mathcal{M}_{K+1}) \geq 2(\delta_{C^2} p_{\min})^{-1}$.

$$2(K^*)^{1/2}\{\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2\}$$

$$\begin{aligned} & \phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(M^*, Y, \Omega)\} \\ & \leq (2K^*)^{1/2}\{\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2\}d(M^*, \mathcal{M}_{K+1}) - \delta_{C^2}p_{\min}d^2(M^*, \mathcal{M}_{K+1}) \\ & \leq -\frac{1}{2}\delta_{C^2}p_{\min}d^2(M^*, \mathcal{M}_{K+1}). \end{aligned} \quad (73)$$

Note that $\phi\{l(\hat{M}^{(K^*)}, Y, \Omega) - l(M^*, Y, \Omega)\} \geq 0$. Thus, the above inequality implies that on the event $d(M^*, \mathcal{M}_K) \geq 4(\delta_{C^2}p_{\min})^{-1}(K^*)^{1/2}\{\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2\}$,

$$\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(\hat{M}^{(K^*)}, Y, \Omega)\} \leq -\frac{1}{2}\delta_{C^2}p_{\min}d^2(M^*, \mathcal{M}_K). \quad (74)$$

Now we proceed to a lower bound for $d(M^*, \mathcal{M}_{K+1})$. Recall the well-known fact that

$\inf_{M \text{ has a rank } K+1} \|M^* - M\|_F^2 = \sum_{l=K+2}^{K^*+1} \sigma_l^2(M^*)$ where $\sigma_1(M^*) \geq \dots \geq \sigma_{K^*+1}(M^*)$ denotes the non-zero singular values of M^* . Thus, $d(M^*, \mathcal{M}_{K+1}) \geq \{\sum_{l=K+2}^{K^*+1} \sigma_l^2(M^*)\}^{1/2} \geq \sigma_{K^*+1}(M^*)$.

Combine this with (74), we have

$$\phi\{l(\hat{M}^{(K)}, Y, \Omega) - l(\hat{M}^{(K^*)}, Y, \Omega)\} \leq -\frac{1}{2}\delta_{C^2}p_{\min}\left\{\sum_{l=K+2}^{K^*+1} \sigma_l^2(M^*)\right\}, \quad (75)$$

if $4(\delta_{C^2}p_{\min})^{-1}(K^*)^{1/2}\{\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2\} \leq \sigma_{K^*+1}(M^*)$ and $K \leq K^* - 1$. We complete the proof by noting that $4(\delta_{C^2}p_{\min})^{-1}(K^*)^{1/2}\{\|Z \circ \Omega\|_2 + 2\delta_{C^2}C^2\|Q\|_2\} \leq \sigma_{K^*+1}(M^*)$ is equivalent to $4(\delta_{C^2}p_{\min})^{-1}K^*R \leq \sigma_{K^*+1}^2(M^*)$. \square

Proof of Lemma 5. According to Lemma 9, for each $0 \leq K \leq K^* - 1$,

$$\begin{aligned} & -2l(\hat{M}^{(K)}, Y, \Omega) + v(n, N, J, K) - \{-2l(\hat{M}^{(K^*)}, Y, \Omega) + v(n, N, J, K^*)\} \\ & \geq \phi^{-1}\delta_{C^2}p_{\min}\left\{\sum_{l=K+2}^{K^*+1} \sigma_l^2(M^*)\right\} - \sum_{l=K+1}^{K^*} u(n, N, J, l), \end{aligned} \quad (76)$$

if $4(\delta_{C^2}p_{\min})^{-1}K^*R \leq \sigma_{K^*}^2(M^*)$. Clearly, right-hand-side of the above inequality is strictly

greater than zero if $u(n, N, J, l) < \phi^{-1} \delta_{C^2} p_{\min} \sigma_{l+1}^2(M^*)$ for all $1 \leq l \leq K^*$. Thus,

$$\begin{aligned}
& \Pr(\hat{K} \geq K^*) \\
& \geq \Pr\left(\max_{1 \leq K \leq K^*} [-2l(\hat{M}^{(K)}, Y, \Omega) + v(n, N, J, K) - \{-2l(\hat{M}^{(K^*)}, Y, \Omega) + v(n, N, J, K^*)\}] > 0\right) \\
& \geq \Pr\left(4(\delta_{C^2} p_{\min})^{-1} K^* R \leq \sigma_{K^*+1}^2(M^*) \text{ and } u(n, N, J, K) < \phi^{-1} \delta_{C^2} p_{\min} \sigma_{K+1}^2(M^*) \text{ for all } 1 \leq K \leq K^*\right)
\end{aligned} \tag{77}$$

The right-hand-side of the above inequality tend to one under the assumptions of the Lemma.

This completes the proof. \square

Proof of Lemma 6. According to Lemma 3, there is a universal constant c such that with probability least $1 - (N + J)^{-1} - (n^*)^{-1}$,

$$\|Z \circ \Omega\|_2 \leq 4(\phi \kappa_{2C^2})^{1/2} (\max_i n_{i\cdot}^*)^{1/2} \vee (\max_j n_{\cdot j}^*)^{1/2} + c\{(\phi/C^2) \vee (\phi \kappa_{2C^2})^{1/2}\} \log n^* \log^{1/2}(N + J). \tag{78}$$

Under the asymptotic regime (8), we have $4(\phi \kappa_{2C^2})^{1/2} = O(1)$, $\max_i n_{i\cdot}^* = O(p_{\max} J)$, $\max_j n_{\cdot j}^* = O(p_{\max} N)$, $c\{(\phi/C^2) \vee (\phi \kappa_{2C^2})^{1/2}\} = O(1)$, and $\log n^* \log^{1/2}(N + J) = O((N \wedge J)^{-1/2} (n^*)^{1/2}) = O(\{p_{\max}(N \vee J)\}^{1/2})$. Thus, the right-hand-side of (78) is of the order $O(\{p_{\max}(N \vee J)\}^{1/2})$ and

$$\|Z \circ \Omega\|_2 = O_p(\{p_{\max}(N \vee J)\}^{1/2}) \tag{79}$$

as $N, J \rightarrow \infty$. Similarly, according to Lemma 2,

$$\|Q\|_2 \leq 4(\max_i n_{i\cdot}^*)^{1/2} \vee (\max_j n_{\cdot j}^*)^{1/2} + c \log^{1/2}(N + J) \tag{80}$$

with probability at least $1 - (N + J)^{-1}$. Under the the asymptotic regime (8), the right-

hand-side of the above inequality is of the order $O(\{p_{\max}(N \vee J)\}^{1/2})$, and thus

$$\|Q\|_2 = O_p(\{p_{\max}(N \vee J)\}^{1/2}). \quad (81)$$

We complete the proof by combining (79), (81), and the definition of R . \square

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723.
- Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221.
- Bandeira, A. S., Van Handel, R., et al. (2016). Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *The Annals of Probability*, 44(4):2479–2506.
- Barrett, P. T., Petrides, K. V., Eysenck, S. B., and Eysenck, H. J. (1998). The Eysenck Personality Questionnaire: An examination of the factorial similarity of P, E, N, and L across 34 countries. *Personality and Individual Differences*, 25:805–819.
- Bartholomew, D. J., Knott, M., and Moustaki, I. (2011). *Latent variable models and factor analysis: A unified approach*. Wiley, Hoboken, NJ.
- Bhaskar, S. A. and Javanmard, A. (2015). 1-bit matrix completion under exact low-rank constraint. In *2015 49th Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE.

- Bock, R. D. and Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: Application of an EM algorithm. *Psychometrika*, 46(4):443–459.
- Buja, A. and Eyuboglu, N. (1992). Remarks on parallel analysis. *Multivariate Behavioral Research*, 27:509–540.
- Cai, T. and Zhou, W.-X. (2013). A max-norm constrained minimization approach to 1-bit matrix completion. *The Journal of Machine Learning Research*, 14:3619–3647.
- Cattell, R. B. (1966). The scree test for the number of factors. *Multivariate Behavioral Research*, 1:245–276.
- Chatterjee, S. et al. (2015). Matrix estimation by universal singular value thresholding. *The Annals of Statistics*, 43(1):177–214.
- Chen, Y., Li, X., and Zhang, S. (2019a). Joint maximum likelihood estimation for high-dimensional exploratory item factor analysis. *psychometrika*, 84(1):124–146.
- Chen, Y., Li, X., and Zhang, S. (2019b). Structured latent factor analysis for large-scale data: Identifiability, estimability, and their implications. *Journal of the American Statistical Association*, pages 1–15.
- Davenport, M. A., Plan, Y., van den Berg, E., and Wootters, M. (2014). 1-bit matrix completion. *Information and Inference*, 3:189–223.
- Dobriban, E. and Owen, A. B. (2019). Deterministic parallel analysis: an improved method for selecting factors and principal components. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 81:163–183.
- Embretson, S. E. and Reise, S. P. (2000). *Item response theory for psychologists*. Lawrence Erlbaum Associates Publishers, Mahwah, NJ.
- Eysenck, S. B., Eysenck, H. J., and Barrett, P. (1985). A revised version of the psychoticism scale. *Personality and Individual Differences*, 6(1):21–29.

- Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, 30(2):179–185.
- Horn, R. A. (1995). Norm bounds for hadamard products and an arithmetic-geometric mean inequality for unitarily invariant norms. *Linear Algebra and Its Applications*, 223:355–361.
- Kaiser, H. F. (1960). The application of electronic computers to factor analysis. *Educational and Psychological Measurement*, 20(1):141–151.
- Little, R. J. and Rubin, D. B. (2019). *Statistical analysis with missing data*. John Wiley & Sons, Hoboken, NJ.
- Ni, R. and Gu, Q. (2016). Optimal statistical and computational rates for one bit matrix completion. In *Artificial Intelligence and Statistics*, pages 426–434.
- Owen, A. B. and Wang, J. (2016). Bi-cross-validation for factor analysis. *Statistical Science*, 31:119–139.
- Reckase, M. (2009). *Multidimensional item response theory*. Springer, New York, NY.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464.
- Skrondal, A. and Rabe-Hesketh, S. (2004). *Generalized latent variable modeling: Multilevel, longitudinal, and structural equation models*. CRC Press, Boca Raton, FL.
- Wedel, M., Böckenholt, U., and Kamakura, W. A. (2003). Factor models for multivariate count data. *Journal of Multivariate Analysis*, 87(2):356–369.
- Zhu, Y., Shen, X., and Ye, C. (2016). Personalized prediction and sparsity pursuit in latent factor models. *Journal of the American Statistical Association*, 111:241–252.