

Testing the Number of Components in Finite Mixture Normal Regression Model with Panel Data

Yu Hao*

Faculty of Business and Economics
The University of Hong Kong
haoyu@hku.hk

Hiroyuki Kasahara

Vancouver School of Economics
The University of British Columbia
hkasahar@mail.ubc.ca

May 15, 2023

Abstract

This paper introduces a likelihood ratio-based test for examining the null hypothesis of an M_0 -component model versus an alternative $(M_0 + 1)$ -component model within the context of normal mixture panel regression. Contrary to the cross-sectional normal mixture, we demonstrate that the first-order derivative of the density function for the variance parameter in the panel normal mixture is linearly independent from its second-order derivative for the mean parameter. However, similar to the cross-sectional normal mixture, the likelihood ratio test statistic for the panel normal mixture remains unbounded. To manage this unboundedness, we employ a penalized maximum likelihood estimator and derive the asymptotic distribution of penalized likelihood ratio test and Expectation-Maximization test statistics using a fourth-order Taylor expansion of the log-likelihood function for reparameterized parameters. A sequential hypothesis testing approach is developed for consistently estimating the number of components. Simulation experiments reveal good finite sample performance of the proposed tests. We apply these tests to estimate the number of production technology types for the finite mixture Cobb-Douglas production function model.

*Address for correspondence: Yu(Jasmine) Hao, Faculty of Business and Economics, The University of Hong Kong. We are very grateful for the comments from Chun Pang Chow, Vadim Marmer, and Kevin Song. We also thank to the IAAE grant at the 2019 IAAE Conference. This research is support by the Natural Science and Engineering Research Council of Canada.

1 Introduction

Finite mixture models offer a natural representation of heterogeneity across a finite number of classes. Due to their flexibility, they have been employed in empirical applications across various fields since the proposal of a two-component normal mixture model by Pearson (1894). In economics, finite mixtures are frequently used to model unobserved individual-specific effects in labor economics, health economics, and industrial organization, among others.¹ Theoretical properties and examples of applications have been discussed by several authors, such as Lindsay (1995), Titterington et al. (1985), and McLachlan and Peel (2004).

The number of components is a crucial parameter in finite mixture models. In economic applications, the number of components often represents the number of unobservable types or abilities. Choosing an arbitrary number of parameters may lead to overestimation or underestimation of the level of heterogeneity. Using too few components may result in biased estimation due to overlooked unobserved heterogeneity, while employing too many components can be computationally costly and ill-behaved because of potential identification problems. Thus, developing a statistical procedure to determine the number of components is essential.

Testing for the number of components in normal mixture regression models has been a long-standing unsolved problem. The regularity conditions of the likelihood ratio test (LRT) for standard asymptotic analysis fail in finite mixture models due to issues such as non-identifiable parameters, the singularity of the Fisher Information matrix, and the true parameter being on the boundary of the parameter space. Numerous papers have been written on the subject of LRT for the number of components (see, e.g., Ghosh and Sen, 1985; Chernoff and Lander, 1995; Lemdani and Pons, 1997; Chen and Chen, 2001, 2003; Chen et al., 2004; Garel, 2001, 2005; Chen et al., 2014), and the asymptotic distribution of the LRT statistic for general finite mixture models has been derived as a function of the Gaussian process (Dacunha-Castelle and Gassiat, 1999; Azaïs et al., 2009; Liu and Shao, 2003; Zhu and Zhang, 2004). However, the key assumptions in these works are violated in cross-sectional normal regression models because normal mixtures possess additional undesirable mathematical properties: (i) the Fisher information for testing is not finite, (ii) the log-likelihood function is unbounded, and (iii) the second derivative of the density function for the mean parameter is linearly dependent on its first derivative for the variance parameter. The asymptotic distribution of the LRT statistics of a cross-sectional univariate finite mixture normal regression model is analyzed by Kasahara and Shimotsu (2015), while its multivariate extension

¹For example, Heckman and Singer (1984) use the finite mixture model to provide an alternative method to account for the unobserved heterogeneity in the analysis of single-spell duration times of unemployed workers. Keane and Wolpin (1997) and Cameron and Heckman (1998) analyze a dynamic model of schooling and occupational choices with unobserved heterogeneous human capital. Likewise, finite mixture models have been applied in health economics. Deb and Trivedi (1997) develop a finite mixture negative binomial count model that accounts for unobserved dispersion of elderly medical care utilization. In industrial organizations, modelling consumer segmentation in marketing such as Kamakura and Russell (1989) and Andrews and Currim (2003) is a venue of application.

is developed by Kasahara and Shimotsu (2019). Amengual et al. (2022) develops a score-type test for a cross-sectional normal mixture model.

This paper develops a likelihood-ratio-based test for determining the number of components in finite mixture normal regression models with panel data, where outcome variables are conditionally independent across periods given latent type within each unit. To the best of our knowledge, it is not known in the existing literature whether the aforementioned problems (i)-(iii) of the cross-sectional normal mixture still arise in the panel normal mixture or not. Furthermore, no likelihood-based test has yet been developed for testing the null hypothesis of an M_0 -component model against an alternative $(M_0 + 1)$ -component model for $M_0 \geq 1$ in the panel normal regression mixture models with conditional independent errors.²

We show that problems related to (i) and (ii) arise, but the higher-order degeneracy of problem (iii) disappears in the panel normal mixture models with conditional independence. Following Chen and Li (2009) and Kasahara and Shimotsu (2015), we consider a penalized likelihood ratio test (PLRT) and an Expectation-Maximization (EM) test to deal with the unboundedness and analyze the asymptotic distribution of the PLRT using a reparameterization orthogonal to the direction in which the Fisher information matrix is singular. The likelihood ratio of an $(M_0 + 1)$ -component model against the M_0 -component model is approximated with local quadratic-form expansion with squares and cross-products of the reparameterized parameters. We demonstrate that the asymptotic null distributions of the penalized likelihood ratio test statistic (PLRTS) and the EM test statistic are characterized by the maximum of M_0 random variables, which we can easily simulate. Building on the PLRT and EM tests, we propose a sequential hypothesis testing approach for consistently estimating the number of components. In simulations, our proposed PLRT and EM tests demonstrate favorable finite sample properties. Moreover, a sequential hypothesis testing approach accurately selects the correct number of components with high frequency, surpassing selection procedures based on the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC).

This paper makes several contributions. First, it analyzes the likelihood-ratio-based test for the number of components in the panel normal regression mixture models with conditional independence. Kasahara and Shimotsu (2015) and Kasahara and Shimotsu (2019) analyze the likelihood-ratio-based tests for the number of components in the cross-sectional univariate normal mixture regression models and the multivariate normal mixture models, respectively. We demonstrate that the asymptotic distribution of PLRT and EM test for the panel normal regression mixture models with conditionally independent errors differs from those of the univariate/multivariate normal mixture models in the aforementioned two papers because the higher-order dependency does not occur when the repeated measurement of outcome variables is available in panel data. Further-

²Kasahara and Shimotsu (2014) develops a procedure to estimate a lower bound on the number of components consistently in finite mixture models in which each component distribution has independent marginals, which includes the panel normal regression mixture models with conditionally independent errors as a special case.

more, we develop a sequential hypothesis testing approach for consistently estimating the number of components.

Second, while it is well known that the log-likelihood function of normal mixture models is unbounded (Hartigan, 1985), it is unknown if the related unboundedness problem arises in the panel data. We show that the likelihood ratio test statistic is unbounded in the panel normal mixture models with conditionally independent errors when the time dimension of panel data is finite. This unboundedness causes over-rejection of the likelihood ratio test. We introduce a penalty function to prevent the likelihood ratio test statistics from being unbounded, where we use computational experiments to determine the data-driven penalty function. We develop an R package `NormalRegPanelMixture` (Hao, 2017) that contains the EM test module and asymptotic distribution simulation module.

Third, our empirical analysis of the number of production technology types using panel data from Japanese and Chilean manufacturing firms provides strong evidence for substantial heterogeneity in production function coefficients across firms within narrowly defined industries. This is an important contribution to the literature on production function estimation, where most existing empirical applications assume the homogeneity of production function coefficients across firms using the standard production function estimation methods developed by Olley and Pakes (1996), Levinsohn and Petrin (2003), and Akerberg et al. (2015). Our empirical finding suggests that it is essential to incorporate unobserved heterogeneity in the production function coefficients across firms in applications (Li and Sasaki, 2017; Doraszelski and Jaumandreu, 2018; Balat et al., 2019; Kasahara et al., 2022).

The EM test approach was introduced by Li et al. (2009) and Chen and Li (2009) to test homogeneity in finite mixture models. Li and Chen (2010) developed an EM test for the null hypothesis of M_0 components applicable to general $M_0 \geq 2$, while Kasahara and Shimotsu (2015) proposed an EM test for normal regression mixture models to test the null of $M_0 \geq 2$. The EM approach has also been applied to test homogeneity in multivariate mixtures (Niu et al., 2011) and subgroup analyses (Shen and He, 2015). More recently, Liu et al. (2018) extended the EM test to mixtures of the general location-scale family distribution, and Kasahara and Shimotsu (2019) developed an EM test for multivariate normal mixture models. Building upon the prior literature, this paper develops an EM test for panel normal regression mixture models with conditionally independent errors.

Identifying and estimating latent group structure in panel data have received attention in recent literature (Kasahara and Shimotsu, 2009; Ando and Bai, 2016; Bonhomme and Manresa, 2015; Lin and Ng, 2012; Lu and Su, 2017; Su et al., 2016). Finite mixture modeling provides a practical, model-based approach to determining unobserved group structures. Choosing the number of groups is often a prerequisite for classifying each individual's group membership. We can estimate the number of groups in panel data regression models by applying our proposed sequential

hypothesis testing approach.

The rest of the paper is organized as follows. In Section 2, we define the finite normal mixture panel regression model. In Section 3, we demonstrate the PLRT for testing the homogeneity of normal mixture panel regression against a two-component model as a precursor to obtaining the general M_0 components test. Section 4 generalizes the result to testing M_0 components against $M_0 + 1$ components. Section 5 introduces the EM test for testing M_0 components against $M_0 + 1$ components. Section 6 derives the asymptotic distribution of the PLRT and EM tests under local alternatives, while Section 7 develops a consistent estimator for the number of components based on sequential hypothesis testing. Section 8 presents the simulated results of the tests. Section 9 provides an empirical application. Let $:=$ denote “equals by definition.” Boldface letters denote vectors or matrices.

2 Heteroskedastic finite mixture panel normal regression model

We consider finite mixture normal regression models with panel data, where the panel length T is fixed and the number of cross-sectional observations n goes to infinity. Define $\mathbf{w} := \{y_t, \mathbf{x}_t, \mathbf{z}_t\}_{t=1}^T$ with $y_t \in \mathbb{R}, \mathbf{x}_t \in \mathbb{R}^q, \mathbf{z}_t \in \mathbb{R}^p$. Given $M \geq 2$, denote the density of a M -component model that represents the conditional density function of $\{y_t\}_{t=1}^T$ given $\{\mathbf{x}_t, \mathbf{z}_t\}_{t=1}^T$ as

$$f_M(\mathbf{w}; \boldsymbol{\vartheta}_M) = \sum_{j=1}^M \alpha_j f(\mathbf{w}; \gamma, \boldsymbol{\theta}_j), \quad (1)$$

where $\boldsymbol{\vartheta}_M = (\boldsymbol{\alpha}^\top, \boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_M^\top, \gamma^\top)^\top \in \Theta_{\boldsymbol{\vartheta}_M}$, $\boldsymbol{\alpha}^\top := (\alpha_1, \dots, \alpha_{M-1})$, $\alpha_M = 1 - \sum_{j=1}^{M-1} \alpha_j$, and

$$f(\mathbf{w}; \gamma, \boldsymbol{\theta}_j) = \prod_{t=1}^T \frac{1}{\sigma_j} \phi\left(\frac{y_t - \mu_j - \mathbf{x}_t^\top \boldsymbol{\beta}_j - \mathbf{z}_t^\top \gamma}{\sigma_j}\right) \quad (2)$$

is the j -th component density function with $\mu_j \in \Theta_\mu \subset \mathbb{R}$, $\sigma_j^2 \in \Theta_\sigma \subset \mathbb{R}_{++}$, $\boldsymbol{\beta}_j \in \Theta_\beta \subset \mathbb{R}^q$, $\gamma \in \Theta_\gamma \subset \mathbb{R}^p$, and $\phi(t) = (2\pi)^{-1/2} \exp(-\frac{t^2}{2})$ is the standard normal probability density function. We collect the component-specific parameters into $\boldsymbol{\theta}_j := (\mu_j, \sigma_j^2, \boldsymbol{\beta}_j^\top)^\top \in \Theta_\theta$ while the regression coefficient γ for a vector \mathbf{z} is assumed to be common across components.

The number of components, denoted by M_0 , is defined as the smallest integer M such that the data density of \mathbf{w} admits the representation (1). Consider a random sample of n with panel length of T independent observations $\{\mathbf{W}_i\}_{i=1}^n$ where $\mathbf{W}_i = \{(Y_{it}, \mathbf{X}_{it}^\top, \mathbf{Z}_{it}^\top)^\top\}_{t=1}^T$ from a true M_0 -component density $f_M(\mathbf{w}; \boldsymbol{\vartheta}_{M_0}^*)$ defined in equation (1) with $\boldsymbol{\vartheta}_{M_0}^* = ((\boldsymbol{\alpha}^*)^\top, (\boldsymbol{\theta}_1^*)^\top, \dots, (\boldsymbol{\theta}_{M_0}^*)^\top, (\gamma^*)^\top)^\top$. The superscript $*$ signifies the true parameter value. Because component distributions can be identified only up to permutation, we assume that $\mu_1^* < \mu_2^* <$

$\dots < \mu_{M_0}^*$ for identification.³

Our goal is to test

$$H_0 : M = M_0 \text{ against } H_A : M = M_0 + 1.$$

3 Likelihood ratio test for $H_0 : M = 1$ against $H_A : M = 2$

We begin by developing the PLRT to test the null hypothesis $H_0 : M = 1$ against the alternative hypothesis $H_1 : M = 2$. Consider a random sample of n with a panel length of T independent observations $\{\mathbf{W}_i\}_{i=1}^n$, where $\mathbf{W}_i = \{(Y_{it}, \mathbf{X}_{it}^\top, \mathbf{Z}_{it}^\top)^\top\}_{t=1}^T$, drawn from a true one-component density $f(\mathbf{w}; \gamma, \boldsymbol{\theta})$ defined in equation (2). Now consider a two-component mixture density function

$$f_2(\mathbf{w}; \boldsymbol{\vartheta}_2) = \alpha f(\mathbf{w}; \gamma, \boldsymbol{\theta}_1) + (1 - \alpha) f(\mathbf{w}; \gamma, \boldsymbol{\theta}_2),$$

where $\boldsymbol{\vartheta}_2 = (\alpha, \boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_2^\top, \gamma^\top)^\top \in \Theta_{\boldsymbol{\vartheta}_2}$, and α is the mixing probability of the first component. The two-component model can generate the true one-component density in two cases: (1) $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}^*$; (2) $\alpha = 0$ or 1 . Consequently, the null hypothesis $H_0 : M = 1$ can be partitioned into two sub-hypotheses: $H_{01} : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$ and $H_{02} : \alpha(1 - \alpha) = 0$. The regularity conditions of the LRTS for a standard asymptotic analysis fail in finite mixture models: under H_{01} , α is not identified, and the Fisher information matrix for the other parameters becomes singular; under H_{02} , α is on the boundary of the parameter space, and either $\boldsymbol{\theta}_1$ or $\boldsymbol{\theta}_2$ is not identified.

As discussed in the introduction, analyzing the asymptotic distribution of the LRTS for the cross-sectional normal mixture is challenging due to its undesirable mathematical properties (cf., Chen and Li, 2009): (i) the Fisher information for testing H_{02} is not finite, (ii) the log-likelihood function is unbounded (Hartigan, 1985), and (iii) the first-order derivative of $f_2(\mathbf{w}; \boldsymbol{\vartheta}_2)$ with respect to σ_j^2 is linearly dependent on its second-order derivative with respect to μ_j . The presence of problems (i)-(iii) in panel normal mixture models with $T \geq 2$ is not well-understood in the existing literature because, to the best of our knowledge, no existing studies have examined them.

Regarding problem (i), we note that the issue of infinite Fisher information for testing H_{02} also arises in the panel normal mixture model. For brevity, let us consider the case without (\mathbf{X}, \mathbf{Z}) . The score for testing $H_{02} : \alpha = 0$ takes the form

$$\left. \frac{\partial f_2(\mathbf{W}; \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)}{\partial \alpha} \right|_{\alpha=0, \mu_2=\mu^*, \sigma_2^2=\sigma^{*2}} = \frac{f(\mathbf{W}; \mu_1, \sigma_1^2)}{f(\mathbf{W}; \mu^*, \sigma^{*2})} - 1,$$

where $f(\mathbf{W}; \mu, \sigma^2) = \prod_{t=1}^T \phi((y_t - \mu)/\sigma)/\sigma$, and $\phi(\cdot)$ is the standard normal density function.

³More generally, we may consider a lexicographical order: $\boldsymbol{\theta}_1^* < \boldsymbol{\theta}_2^* < \dots < \boldsymbol{\theta}_{M_0}^*$.

When $\sigma_1^2 > 2\sigma^{*2}$, $\mathbb{E}[\{f(\mathbf{W}; \mu_1, \sigma_1^2)/f(\mathbf{W}; \mu^*, \sigma^{*2}) - 1\}^2] = \infty$. For more details, please refer to Proposition 5. Because the infinite Fisher information causes difficulty in deriving the asymptotic distribution under H_{02} , this paper focuses on testing H_{01} . We define $\Upsilon_1^* := \{(\alpha, \gamma, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \in \Theta_{\vartheta_2} : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}^* \text{ and } \gamma = \gamma^*\}$, which is the subspace of Θ_{ϑ_2} that corresponds to H_{01} . Note that because of our focus on H_{01} , our test may not have power against the local alternatives with $\alpha_n \rightarrow 0$. We analyze the asymptotic distribution of the PLRTS under the contiguous local alternatives in Section 6.

Related to problem (ii), the LRTS in normal mixture models with panel data becomes unbounded as the sample size n goes to ∞ . Define the likelihood ratio statistic with respect to the true parameter under H_0 as:

$$LR_n^*(\boldsymbol{\vartheta}_2) := 2 \left\{ \sum_{i=1}^n \log f_2(\mathbf{W}_i; \boldsymbol{\vartheta}_2) - \sum_{i=1}^n \log f(\mathbf{W}_i; \gamma^*, \boldsymbol{\theta}^*) \right\},$$

where f_2 is the density of the two-component finite mixture distribution in (1) with $M = 2$, and $((\gamma^*)^\top, (\boldsymbol{\theta}^*)^\top)^\top$ is the true parameter value under H_0 . Let $\tilde{\boldsymbol{\vartheta}}_{2,n}$ be the maximum likelihood estimator for the two-component model, i.e., $\tilde{\boldsymbol{\vartheta}}_{2,n} = \arg \max_{\boldsymbol{\vartheta}_2 \in \Theta_{\boldsymbol{\vartheta}_2}} LR_n^*(\boldsymbol{\vartheta}_2)$.

Proposition 1. *Suppose the true model is described by the one-component model with $(\gamma, \boldsymbol{\theta}) = (\gamma^*, \boldsymbol{\theta}^*)$. Then, for any positive constant $M > 0$, $\Pr \left(LR_n^*(\tilde{\boldsymbol{\vartheta}}_{2,n}) \leq M \right) \rightarrow 0$ as $n \rightarrow \infty$.*

To deal with unboundedness, we consider a maximum penalized likelihood estimator (PMLE) as in Chen and Tan (2009) using the following penalty function:

$$\tilde{p}_n(\boldsymbol{\vartheta}_M) := \sum_{j=1}^M p_n(\sigma_j^2) \quad \text{with} \quad p_n(\sigma_j^2) := -a_n \{ \sigma_0^2 / \sigma_j^2 + \log(\sigma_j^2 / \sigma_0^2) - 1 \}. \quad (3)$$

This penalty function circumvents the problem of unbounded log likelihood by preventing a variance parameter estimate from nearing zero. The parameter a_n is selected such that the penalty's impact becomes asymptotically negligible for the distribution of the PMLE. Refer to the conditions C1-C3 in the proof of Proposition 6.

Let

$$\hat{\boldsymbol{\vartheta}}_2 = \arg \max_{\boldsymbol{\vartheta}_2 \in \Theta_{\boldsymbol{\vartheta}_2}} \sum_{i=1}^n \log f_2(\mathbf{W}_i; \boldsymbol{\vartheta}_2) + \tilde{p}_n(\boldsymbol{\vartheta}_2)$$

denote the Penalized Maximum Likelihood Estimator (PMLE, hereafter) under the two-component model. Define a set of parameter values for the two-component density that generates the true one-component density by $\Theta_2^* := \{(\alpha, \gamma, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \in \Theta_{\vartheta_2} : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}^* \text{ and } \gamma = \gamma^*; \alpha = 1 \text{ and } \boldsymbol{\theta}_1 = \boldsymbol{\theta}^*; \alpha = 0 \text{ and } \boldsymbol{\theta}_2 = \boldsymbol{\theta}^*\}$. $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ are component-specific parameters, while γ is a parameter vector common across components. The following proposition establishes the consistency

of the PMLE.

Assumption 1. (a) \mathbf{X} and \mathbf{Z} have finite second moments, and $\Pr(\mathbf{X}^\top \boldsymbol{\beta} + \mathbf{Z}^\top \boldsymbol{\gamma} \neq \mathbf{X}^\top \boldsymbol{\beta}^* + \mathbf{Z}^\top \boldsymbol{\gamma}^*) > 0$ for $(\boldsymbol{\beta}^\top, \boldsymbol{\gamma}^\top)^\top \neq ((\boldsymbol{\beta}^*)^\top, (\boldsymbol{\gamma}^*)^\top)^\top$. (b) $a_n > 0$ and $a_n = o(n^{1/4})$ in the penalty function (3).

Proposition 2. Suppose that Assumption 1 holds. Then under the null hypothesis $H_0 : M_0 = 1$, $\inf_{\boldsymbol{\vartheta}_2 \in \Theta_2^*} |\hat{\boldsymbol{\vartheta}}_2 - \boldsymbol{\vartheta}_2| \rightarrow_p 0$.

It should be noted that $f_2(\mathbf{w}; \boldsymbol{\vartheta}_2^*) = f(\mathbf{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)$ for any $\boldsymbol{\vartheta}_2^* \in \Theta_2^*$. Consequently, Proposition 2 suggests that the PMLE $\hat{\boldsymbol{\vartheta}}_2$ converges in probability to a set of parameters for which the true density function $f(\mathbf{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)$ emerges within the space of two-component density functions.

For problem (iii), we show that in normal mixture models with panel data, the first-order derivative of $f_2(\mathbf{w}; \boldsymbol{\vartheta}_2)$ with respect to σ_j^2 is *not* linearly dependent with its second-order derivative with respect to μ_j (See Proposition 3(c)). Consequently, the panel mixture model (1) with the component density function (2) is strongly identifiable, and the best rate of convergence for estimating the mixing distribution is $n^{-1/4}$ when the number of components is unknown (c.f., Chen, 1995). See Proposition 4(a). In contrast, the strong identifiability does not hold for the cross-sectional normal mixture, and its convergence rate becomes as slow as $n^{-1/8}$ when the number of components is over-specified (c.f., Kasahara and Shimotsu, 2015).

As in any finite mixture models, however, the standard asymptotic analysis breaks down in testing $H_{01} : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}^*$ because α is not identified under H_{01} ; in addition, the first-order derivative at the true value $\boldsymbol{\vartheta}_2^* = (\alpha, (\boldsymbol{\theta}^*)^\top, (\boldsymbol{\theta}^*)^\top, (\boldsymbol{\gamma}^*)^\top)^\top$ are linear dependent as

$$\nabla_{\boldsymbol{\theta}_1} \log f_2(\mathbf{w}; \boldsymbol{\vartheta}_2^*) = \frac{\alpha}{1 - \alpha} \nabla_{\boldsymbol{\theta}_2} \log f_2(\mathbf{w}; \boldsymbol{\vartheta}_2^*). \quad (4)$$

To deal with this linear dependency, we analyze the asymptotic distribution of LRTS by developing a higher-order approximation for the log-likelihood function.

To extract the direction of Fisher Information matrix singularity, we adapt the reparameterization approach by Kasahara and Shimotsu (2012) and consider the following one-to-one reparameterization of $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ given α :

$$\begin{pmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\nu} \end{pmatrix} := \begin{pmatrix} \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \\ \alpha \boldsymbol{\theta}_1 + (1 - \alpha) \boldsymbol{\theta}_2 \end{pmatrix} \text{ so that } \begin{pmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{pmatrix} = \begin{pmatrix} \boldsymbol{\nu} + (1 - \alpha) \boldsymbol{\lambda} \\ \boldsymbol{\nu} - \alpha \boldsymbol{\lambda} \end{pmatrix}, \quad (5)$$

where $\boldsymbol{\nu}$ and $\boldsymbol{\lambda}$ are both $(q + 2) \times 1$ reparameterized parameter vectors with $\boldsymbol{\nu} = (\nu_\mu, \nu_\sigma, \boldsymbol{\nu}_\beta)^\top$ and $\boldsymbol{\lambda} = (\lambda_\mu, \lambda_\sigma, (\boldsymbol{\lambda}_\beta)^\top)^\top = (\mu_1 - \mu_2, \sigma_1^2 - \sigma_2^2, (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2)^\top)^\top$. We also write $\boldsymbol{\theta}$ and $\boldsymbol{\lambda}$ as $\boldsymbol{\theta} = (\theta_1, \theta_2, \theta_3, \dots, \theta_{q+2})^\top := (\mu, \sigma^2, \beta_1, \dots, \beta_q)^\top$ and $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_{q+2})^\top := (\lambda_\mu, \lambda_\sigma, \lambda_{\beta_1}, \dots, \lambda_{\beta_q})^\top$.

This reparameterization is essential for analyzing the asymptotic distribution of the PLRTS in light of the linear dependency in (4). The reparameterized parameter $\boldsymbol{\lambda}$ captures a deviation from the one component model, where its first-order derivatives of the log density are identically equal

to zero under $H_0 : M = 1$. Consequently, this reparameterization facilitates the derivation of an approximate quadratic-form criterion function, which is based on the fourth-order Taylor series approximation of the log-likelihood function, in order to characterize the asymptotic distribution of the LRTS.

Define the space for reparameterized parameters as

$$\boldsymbol{\psi} := (\boldsymbol{\gamma}^\top, \boldsymbol{\nu}^\top, \boldsymbol{\lambda}^\top)^\top \in \Theta_\psi,$$

where $\Theta_\psi = \{\boldsymbol{\psi} : \boldsymbol{\gamma} \in \Theta_\gamma, \boldsymbol{\nu} + (1 - \alpha)\boldsymbol{\lambda} \in \Theta_\theta, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda} \in \Theta_\theta\}$. Under the null hypothesis $H_{01} : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}^*$, we have $\boldsymbol{\lambda} = (0, \dots, 0)^\top$ and $\boldsymbol{\nu} = \boldsymbol{\theta}^*$. We rewrite the reparameterized parameters under null hypothesis to be $(\boldsymbol{\psi}^*)^\top = ((\boldsymbol{\gamma}^*)^\top, (\boldsymbol{\theta}^*)^\top, 0, \dots, 0)^\top$. Under the reparameterized parameter space, the density function and its logarithm are expressed as

$$\begin{aligned} g(\boldsymbol{w}; \boldsymbol{\psi}, \alpha) &= \alpha f(\boldsymbol{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} + (1 - \alpha)\boldsymbol{\lambda}) + (1 - \alpha)f(\boldsymbol{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda}) \quad \text{and} \\ l(\boldsymbol{w}; \boldsymbol{\psi}, \alpha) &= \log g(\boldsymbol{w}; \boldsymbol{\psi}, \alpha). \end{aligned} \quad (6)$$

Write $\boldsymbol{\psi}$ as $\boldsymbol{\psi} = (\boldsymbol{\eta}^\top, \boldsymbol{\lambda}^\top)^\top$ with $\boldsymbol{\eta} = (\boldsymbol{\gamma}^\top, \boldsymbol{\nu}^\top)^\top$, where $\boldsymbol{\eta}^* = ((\boldsymbol{\gamma}^*)^\top, (\boldsymbol{\nu}^*)^\top)^\top$ and $\boldsymbol{\lambda}^* = \mathbf{0}$. Denote the parameter space of $\boldsymbol{\eta}$ and $\boldsymbol{\lambda}$ by $\Theta_\eta \subset \mathbb{R}^{p+q+2}$ and $\Theta_\lambda \subset \mathbb{R}^{q+2}$, respectively.

Under this reparameterization, the first-order derivatives of the reparameterized log density with respect to the reparameterized parameters $\boldsymbol{\eta}$ is identical to those under the one-component model, and the first-order derivative with respect to $\boldsymbol{\lambda}$ is a zero vector:

$$\nabla_{\boldsymbol{\eta}^\top} l(\boldsymbol{w}; \boldsymbol{\psi}^*, \alpha) = \frac{\nabla_{(\boldsymbol{\gamma}^\top, \boldsymbol{\theta}^\top)^\top} f(\boldsymbol{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)}{f(\boldsymbol{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)} \quad \text{and} \quad \nabla_{\boldsymbol{\lambda}} l(\boldsymbol{w}; \boldsymbol{\psi}^*, \alpha) = \mathbf{0}. \quad (7)$$

With $\nabla_{\boldsymbol{\lambda}} l(\boldsymbol{w}; \boldsymbol{\psi}^*, \alpha) = \mathbf{0}$, the Fisher information matrix is singular, and the standard quadratic approximation fails. Consequently, the information on $\boldsymbol{\lambda}$ is provided by the second-order derivative of $l(\boldsymbol{w}; \boldsymbol{\psi}, \alpha)$ with respect to $\boldsymbol{\lambda}$. We use second-order derivative with respect to $\boldsymbol{\lambda}$ to identify $\boldsymbol{\lambda}$:

$$\nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}^\top} l(\boldsymbol{w}; \boldsymbol{\psi}^*, \alpha) = \alpha(1 - \alpha) \frac{\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top} f(\boldsymbol{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)}{f(\boldsymbol{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)}. \quad (8)$$

When α is bounded away from 0 and 1, the elements of $\nabla_{\boldsymbol{\lambda}\boldsymbol{\lambda}^\top} l(\boldsymbol{W}; \boldsymbol{\psi}^*, \alpha)$ are mean-zero random variables.

Note that, unlike the cross-sectional models analyzed by Kasahara and Shimotsu (2015), there exists no collinearity between these first and second-order derivatives for the panel models. This distinction is indeed important as it highlights the differences in the asymptotic distribution of the LRTS for the panel models compared to the cross-sectional models. The absence of collinearity between the first and second-order derivatives in panel models leads to different convergence

rates and asymptotic properties.

Let f^* and ∇f^* denote $f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)$ and $\nabla f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)$. Define the vector $\mathbf{s}(\mathbf{W})$ as

$$\mathbf{s}(\mathbf{W}) = \begin{pmatrix} \mathbf{s}_\eta(\mathbf{W}) \\ \mathbf{s}_{\lambda\lambda}(\mathbf{W}) \end{pmatrix}, \quad \text{where } \mathbf{s}_\eta(\mathbf{W}) := \frac{\nabla_{(\boldsymbol{\gamma}^\top, \boldsymbol{\theta}^\top)^\top} f^*}{f^*} \quad \text{and} \quad \mathbf{s}_{\lambda\lambda}(\mathbf{W}) := \frac{\tilde{\nabla}_{\boldsymbol{\theta}\boldsymbol{\theta}^\top} f^*}{f^*}. \quad (9)$$

The term $\tilde{\nabla}_{\boldsymbol{\theta}\boldsymbol{\theta}^\top} f^*$ denotes the second-order derivatives of the density function f^* with respect to the parameters $\boldsymbol{\theta}$. Coefficients c_{jk} are employed to adjust the scaling of these second-order derivatives. The function $\mathbf{s}(\mathbf{w})$ comprises the second-order derivatives of the log-likelihood function with respect to the reparameterized parameter $\boldsymbol{\lambda}$. This function, $\mathbf{s}_{\lambda\lambda}(\mathbf{w})$, serves as a score function for identifying $\boldsymbol{\lambda}$. Consequently, $\mathbf{s}(\mathbf{w})$ is referred to as a score function. An explicit expression for the score function $\mathbf{s}(\mathbf{w})$ can be derived using Hermite polynomials, as elaborated in Appendix B.2.

Collect the relevant normalized reparameterized parameters and define $\mathbf{t}(\boldsymbol{\psi}, \alpha)$ as

$$\mathbf{t}(\boldsymbol{\psi}, \alpha) = \begin{pmatrix} \mathbf{t}_\eta \\ \mathbf{t}_\lambda(\boldsymbol{\lambda}, \alpha) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\eta} - \boldsymbol{\eta}^* \\ \alpha(1 - \alpha)\mathbf{v}(\boldsymbol{\lambda}) \end{pmatrix}, \quad (10)$$

where $\mathbf{v}(\boldsymbol{\lambda})$ is a vector of unique elements of $\boldsymbol{\lambda}\boldsymbol{\lambda}^\top$ given by

$$\mathbf{v}(\boldsymbol{\lambda}) = (\lambda_1\lambda_1, \dots, \lambda_{q+2}\lambda_{q+2}, \lambda_1\lambda_2, \dots, \lambda_{q+1}\lambda_{q+2})^\top \quad (11)$$

of which length is $q_\lambda := (q+2)(q+3)/2$.

Let $L_n(\boldsymbol{\psi}, \alpha) := \sum_{i=1}^n l(\mathbf{W}_i; \boldsymbol{\psi}^*, \alpha)$ be the reparameterized log-likelihood function and define the normalized score vector

$$\mathbf{S}_n := n^{-1/2} \sum_{i=1}^n \mathbf{s}(\mathbf{W}_i).$$

Then, taking the fourth order Taylor expansion of $L_n(\boldsymbol{\psi}, \alpha)$ around $(\boldsymbol{\psi}^*, \alpha)$, we may write $2\{L_n(\boldsymbol{\psi}, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha)\}$ as a quadratic function of $\sqrt{n}\mathbf{t}(\boldsymbol{\psi}, \alpha)$ as

$$2\{L_n(\boldsymbol{\psi}, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha)\} = 2(\sqrt{n}\mathbf{t}(\boldsymbol{\psi}, \alpha))^\top \mathbf{S}_n - (\sqrt{n}\mathbf{t}(\boldsymbol{\psi}, \alpha))^\top \mathcal{I}_n(\sqrt{n}\mathbf{t}(\boldsymbol{\psi}, \alpha)) + R_n(\boldsymbol{\psi}, \alpha) \quad (12)$$

$$= \mathbf{G}_n^\top \mathcal{I}_n \mathbf{G}_n - [\sqrt{n}\mathbf{t}(\boldsymbol{\psi}, \alpha) - \mathbf{G}_n]^\top \mathcal{I}_n [\sqrt{n}\mathbf{t}(\boldsymbol{\psi}, \alpha) - \mathbf{G}_n] + R_n(\boldsymbol{\psi}, \alpha), \quad (13)$$

where \mathcal{I}_n is the negative of the sample Hessian defined in the proof of Proposition 3 while $\mathbf{G}_n := \mathcal{I}_n^{-1} \mathbf{S}_n$. Let $\mathcal{I} = \mathbb{E}[\mathbf{s}(\mathbf{W})\mathbf{s}(\mathbf{W})^\top]$.

Assumption 2. (a) \mathbf{X} and \mathbf{Z} have finite 8-th moments. (b) $\mathbb{E}[\mathbf{U}\mathbf{U}^\top]$ is non-singular, where $\mathbf{U} = [1, \mathbf{X}^\top, \mathbf{Z}^\top]^\top$.

Proposition 3. *Suppose that assumption 1 and 2 hold. Then, under $H_0 : M = 1$, for $\alpha \in (0, 1)$, (a) for any $\delta > 0$, $\limsup_{n \rightarrow \infty} \Pr(\sup_{\psi \in \Theta_\psi: \|\psi - \psi^*\| \leq \kappa} |R_n(\psi, \alpha)| > \delta(1 + \|nt(\psi, \alpha)\|^2)) \rightarrow 0$ as $\kappa \rightarrow 0$, (b) $\mathcal{S}_n \xrightarrow{d} \mathcal{S} \sim N(0, \mathcal{I})$, and (c) $\mathcal{I}_n \xrightarrow{p} \mathcal{I}$, where \mathcal{I} is finite and non-singular.*

The non-singularity of \mathcal{I} in Proposition 3(c) highlights the difference between the cross-sectional normal mixture and the panel data normal mixture models. In particular, as shown in equation (75) in Appendix B.2, the first-order derivative of $f_2(\mathbf{w}; \boldsymbol{\vartheta}_2)$ with respect to σ_j^2 is linearly independent of its second-order derivative with respect to μ_j when $T \geq 2$, ensuring that the high-order degeneracy of problem (iii) does not arise. Intuitively, the availability of repeated observations within each individual unit provides better identification, even for over-parameterized models, and reduces the degree of higher-order degeneracy.

The set of feasible values of $\sqrt{n}t(\psi, \alpha)$ is given by the shifted and re-scaled parameter space for $(\boldsymbol{\eta}, v(\boldsymbol{\lambda}))$ defined as $\Lambda_n := \sqrt{n}(\Theta_\eta - \eta^*) \times \sqrt{n}\alpha(1 - \alpha)v(\Theta_\lambda)$, where $v(A) := \{t \in \mathbb{R}^{q_\lambda} : t = v(\lambda) \text{ for some } \lambda \in A \subset \mathbb{R}^{q+2}\}$. Because Λ_n/\sqrt{n} is locally approximated by a cone $\Lambda := \mathbb{R}^{p+q+2} \times v(\mathbb{R}^{q+2})$, we may apply Lemma 2 of Andrews (1999) to approximate the distribution of the supremum of the right-hand side of (13) as

$$\max_{\psi \in \Theta_\psi} 2\{L_n(\psi, \alpha) - L_n(\psi^*, \alpha)\} \xrightarrow{d} \mathbf{G}^\top \mathcal{I} \mathbf{G} - \inf_{\mathbf{t} \in \Lambda} (\mathbf{t} - \mathbf{G})' \mathcal{I} (\mathbf{t} - \mathbf{G}),$$

where $\mathbf{G} = \mathcal{I}^{-1} \mathcal{S} \sim N(0, \mathcal{I}^{-1})$. This allows us to characterize the asymptotic distribution of the LRTS.

For each $\alpha \in (0, 1)$, define the reparameterized PMLE by

$$\hat{\boldsymbol{\psi}} = \arg \max_{\boldsymbol{\psi} \in \Theta_\psi} L_n(\boldsymbol{\psi}, \alpha) + \sum_{j=1}^2 p_n(\sigma_j^2(\boldsymbol{\psi}, \alpha)) \quad (14)$$

with $\hat{\boldsymbol{\psi}} := (\hat{\boldsymbol{\gamma}}^\top, \hat{\boldsymbol{\nu}}^\top, \hat{\boldsymbol{\lambda}}^\top)^\top$, where Θ_ψ is defined as the space of $\boldsymbol{\psi}$ so that the $\boldsymbol{\vartheta}_2$ implied is in Θ_ϑ and $\sigma_j^2(\boldsymbol{\psi}, \alpha)$ is the value of σ_j implied by the value of $\boldsymbol{\psi}$ and α (e.g., $\sigma_1^2(\boldsymbol{\psi}, \alpha) = \nu_\sigma + (1 - \alpha)\lambda_\sigma$).

Let $(\hat{\boldsymbol{\gamma}}_0, \hat{\boldsymbol{\theta}}_0)$ be the one-component MLE that maximizes the one-component likelihood function $L_{0,n}(\boldsymbol{\gamma}, \boldsymbol{\theta}) := \sum_{i=1}^n \log f(\mathbf{W}_i; \boldsymbol{\gamma}, \boldsymbol{\theta})$. Define the LRTS and the PLRTS of testing H_{01} , respectively, with a small positivity constant ϵ on α as

$$LR_n := \max_{\alpha \in [\epsilon, 1 - \epsilon]} 2\{L_n(\hat{\boldsymbol{\psi}}, \alpha) - L_{0,n}(\hat{\boldsymbol{\gamma}}_0, \hat{\boldsymbol{\theta}}_0)\} \text{ and } PLR_n := LR_n + \sum_{j=1}^2 p_n(\sigma_j^2(\hat{\boldsymbol{\psi}}, \alpha)). \quad (15)$$

The hard bound is imposed on the values of α in order to avoid an issue of the infinite Fisher information for testing H_{02} . However, the LRTS may have a reduced power if the true value of α does not satisfy the constraint $[\epsilon, 1 - \epsilon]$ given an ad hoc constant $\epsilon > 0$. For this reason, we also

develop the EM-test in Section 5 which does not impose a direct constraint on the value of α .

With $s(\mathbf{W})$ in (9), partition $\mathcal{I} = \mathbb{E}[s(\mathbf{W})s(\mathbf{W})^\top]$ and define

$$\mathcal{I} = \begin{pmatrix} \mathcal{I}_\eta & \mathcal{I}_{\eta\lambda} \\ \mathcal{I}_{\lambda\eta} & \mathcal{I}_{\lambda\lambda} \end{pmatrix}, \quad \mathcal{I}_\eta = \mathbb{E}[s_\eta(\mathbf{W})s_\eta(\mathbf{W})^\top], \quad \mathcal{I}_{\lambda\eta} = \mathbb{E}[s_{\lambda\lambda}(\mathbf{W})s_\eta(\mathbf{W})^\top], \quad \mathcal{I}_{\eta\lambda} = \mathcal{I}_{\lambda\eta}^\top, \\ \mathcal{I}_{\lambda\lambda} = \mathbb{E}[s_{\lambda\lambda}(\mathbf{W})s_{\lambda\lambda}(\mathbf{W})^\top], \quad \mathcal{I}_{\lambda,\eta} = \mathcal{I}_{\lambda\lambda} - \mathcal{I}_{\lambda\eta}\mathcal{I}_\eta^{-1}\mathcal{I}_{\eta\lambda}, \quad \text{and} \quad \mathbf{G}_{\lambda,\eta} := (\mathcal{I}_{\lambda,\eta})^{-1}\mathbf{S}_{\lambda,\eta},$$

where $\mathbf{S}_{\lambda,\eta} \sim N(0, \mathcal{I}_{\lambda,\eta})$. Define a set that characterizes the feasible values of $\sqrt{n}\mathbf{t}_\lambda(\lambda, \alpha)$ when $n \rightarrow \infty$ by the cone

$$\Lambda_\lambda = \left\{ \sqrt{n}\alpha(1 - \alpha)v(\lambda) : \lambda \in \Theta_\lambda \right\}.$$

Define $\hat{\mathbf{t}}_\lambda$ by

$$r_\lambda(\hat{\mathbf{t}}_\lambda) = \inf_{\mathbf{t}_\lambda \in \Lambda_\lambda} r_\lambda(\mathbf{t}_\lambda), \quad r_\lambda(\mathbf{t}_\lambda) := (\mathbf{t}_\lambda - \mathbf{G}_{\lambda,\eta})^\top \mathcal{I}_{\lambda,\eta} (\mathbf{t}_\lambda - \mathbf{G}_{\lambda,\eta}), \quad (16)$$

where $\hat{\mathbf{t}}_\lambda$ is a projection of a random Gaussian random variable \mathbf{G}_λ on a cone Λ_λ .

The following proposition establishes the asymptotic distribution of LRTS or PLRTS under the null hypothesis $H_0 : M = 1$.

Proposition 4. *Suppose that assumptions 1 and 2 hold. Under the null hypothesis $H_0 : M_0 = 1$, (a) $\mathbf{t}(\hat{\psi}, \alpha) = O_p(n^{-1/2})$ for any $\alpha \in (0, 1)$, (b) $LR_n \xrightarrow{d} (\hat{\mathbf{t}}_\lambda)^\top \mathcal{I}_{\lambda,\eta} \hat{\mathbf{t}}_\lambda$ and $PLR_n \xrightarrow{d} (\hat{\mathbf{t}}_\lambda)^\top \mathcal{I}_{\lambda,\eta} \hat{\mathbf{t}}_\lambda + \text{plim}_{n \rightarrow \infty} \sum_{j=1}^2 p_n(\sigma_j^2(\hat{\psi}, \alpha))$.*

Proposition 4(a) implies that $\hat{\theta}_j - \theta^* = O_p(n^{-1/4})$ for $j = 1, 2$. The $n^{1/4}$ convergence rate is a consequence of the linear dependency in (4), where the identification of the parameter θ relies on the fourth-order Taylor approximation of the log-likelihood function. This rate is also the best convergence rate for an over-parameterized mixture under the strong identifiability condition (Chen, 1995). When we choose the penalty function so that $\sum_{j=1}^2 p_n(\sigma_j^2(\hat{\psi}, \alpha)) = o_p(1)$ under the null hypothesis of $M = 1$, PLR_n has the same asymptotic null distribution as that of LR_n .

4 Likelihood ratio test for $H_0 : M = M_0$ against $H_A : M = M_0 + 1$

In this section, we build upon the analysis from the previous section and derive the asymptotic distribution of the PLRTS for testing the null hypothesis of M_0 components against an alternative of $(M_0 + 1)$ components, where $M_0 \geq 2$.

Consider a random sample of n with a panel length of T independent observations $\{\mathbf{W}_i\}_{i=1}^n$, where $\mathbf{W}_i = \{(Y_{it}, \mathbf{X}_{it}^\top, \mathbf{Z}_{it}^\top)^\top\}_{t=1}^T$ from an M_0 -component density $f_{M_0}(\mathbf{w}; \boldsymbol{\vartheta}_{M_0})$ defined in equa-

tion (17):

$$f_{M_0}(\mathbf{w}; \boldsymbol{\vartheta}_{M_0}^*) = \sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_j^*), \quad (17)$$

where $\boldsymbol{\vartheta}_{M_0}^* = (\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \dots, \boldsymbol{\theta}_{M_0}^*, \alpha_1^*, \dots, \alpha_{M_0-1}^*, \boldsymbol{\gamma}^*) \in \Theta_{\boldsymbol{\vartheta}_{M_0}}$ and $\alpha_{M_0}^* = 1 - \sum_{j=1}^{M_0-1} \alpha_j^*$.

Let the density of the $(M_0 + 1)$ -component model be defined by:

$$f_{M_0+1}(\mathbf{w}; \boldsymbol{\vartheta}_{M_0+1}) = \sum_{j=1}^{M_0+1} \alpha_j f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\theta}_j), \quad (18)$$

where $\boldsymbol{\vartheta}_{M_0+1} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_{M_0+1}, \alpha_1, \dots, \alpha_{M_0}, \boldsymbol{\gamma}) \in \Theta_{\boldsymbol{\vartheta}_{M_0+1}}$ as defined in (17). We assume $\mu_1^* < \mu_2^*, \dots, < \mu_{M_0}^*$ in the true parameters for identification.

The $(M_0 + 1)$ -component model (18) gives rise to the true density (17) in two different cases: (i) two components have the same mixing parameter, and (ii) one component has zero mixing proportion. Accordingly, we partition the null hypothesis of $H_0 : M = M_0$ into two as $H_0 = H_{01} \cup H_{02}$, with $H_{01} : \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1} = \boldsymbol{\theta}_h^*$ for some $h = 1, \dots, M_0$, and $H_{02} : \alpha_h = 0$ for some $h = 1, \dots, M_0 + 1$.

We first analyze the infinite Fisher information problem for testing H_{02} . Partition H_{02} as $H_{02} = \cup_{h=1}^{M_0} H_{0,2h}$, where $H_{0,2h} : \alpha_h = 0$. Define the subset of $\Theta_{\boldsymbol{\vartheta}_{M_0+1}}$ corresponding to $H_{0,2h}$ as

$$\begin{aligned} \Upsilon_{2h}^* = \{ & \boldsymbol{\vartheta}_{M_0+1} \in \Theta_{\boldsymbol{\vartheta}_{M_0+1}} : \alpha_h = 0; (\alpha_j, \mu_j, \sigma_j) = (\alpha_j^*, \mu_j^*, \sigma_j^*) \text{ for } j < h; \\ & (\alpha_j, \mu_j, \sigma_j) = (\alpha_{j-1}^*, \mu_{j-1}^*, \sigma_{j-1}^*) \text{ for } j > h\}. \end{aligned}$$

The score for testing $H_{0,2h} : \alpha_h = 0$ takes the form $\nabla_{\alpha_h} \log f_{M_0+1}(\mathbf{W}_i, \boldsymbol{\vartheta}_{M_0+1}) = [f(\mathbf{W}_i; \mu_h, \sigma_h^2) - f(\mathbf{W}_i; \mu_{M_0}^*, \sigma_{M_0}^{2*})] / f_{M_0}(\mathbf{W}_i, \boldsymbol{\vartheta}_{M_0}^*)$. Because (μ_h, σ_h^2) is not identified when $\alpha_h = 0$, the Fisher information matrix of the LRTS for testing $H_{0,2h} : \alpha_h = 0$ depends on the supremum of the variance of $\nabla_{\alpha_h} \log f_{M_0+1}(\mathbf{W}_i; \boldsymbol{\vartheta}_{M_0+1})$ over $\boldsymbol{\vartheta}_{M_0+1} \in \Upsilon_{2h}^*$. The Fisher information is infinite unless there is an a priori restriction on the values of σ_j^2 .

Proposition 5. $\sup_{\boldsymbol{\vartheta}_{M_0+1} \in \Upsilon_{2h}^*} \mathbb{E}[\{\nabla_{\alpha_h} \log f_{M_0+1}(\mathbf{W}_i, \boldsymbol{\vartheta}_{M_0+1})\}^2] < \infty$ if and only if $\max\{\sigma^2 : \sigma \in \Theta_\sigma\} < 2 \max\{\sigma_1^{2*}, \dots, \sigma_{M_0}^{2*}\}$.

Since the restriction on the values of σ_j^2 in Proposition 5 is difficult to justify and not easy to enforce in practice, we focus on testing H_{01} .

Partition H_{01} as $H_{01} = \cup_{h=1}^{M_0} H_{0,1h}$, where $H_{0,1h} : \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1}$ with $\mu_1 < \dots < \mu_h = \mu_{h+1} < \dots < \mu_{M_0+1}$. We impose these inequality constraints on μ_j for component identification. There are M_0 ways to describe the M_0 component null model in the space of $(M_0 + 1)$ component models, each way corresponding to the null hypothesis of $H_{0,1h} : \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1}$ for $h = 1, 2, \dots, M_0$. Testing $H_{0,1h} : \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1}$ in the M_0 -component null models is similar to testing $H_{01} : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$ in the one

component null model in section 3.

Define the subset of $\Theta_{\vartheta_{M_0+1}}$ corresponding to $H_{0,1h}$ as:

$$\begin{aligned} \Upsilon_{1h}^* := & \left\{ \vartheta_{M_0+1} \in \Theta_{\vartheta_{M_0+1}} : \alpha_h + \alpha_{h+1} = \alpha_h^* \text{ and } \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1} = \boldsymbol{\theta}_h^*; \gamma = \gamma^*; \alpha_j = \alpha_j^* \right. \\ & \left. \text{and } \boldsymbol{\theta}_j = \boldsymbol{\theta}_j^* \text{ for } 1 \leq j < h; \alpha_j = \alpha_{j-1}^* \text{ and } \boldsymbol{\theta}_j = \boldsymbol{\theta}_{j-1}^* \text{ for } h+1 \leq j \leq M_0+1 \right\} \end{aligned} \quad (19)$$

for $h = 1, \dots, M_0$. The set $\Upsilon_1^* := \cup_{h=1}^{M_0} \Upsilon_{1h}^*$ corresponds to $H_{01} = \cup_{h=1}^{M_0} H_{0,1h}$.

Suppose the null hypothesis of $M = M_0$ holds with the true density (17). Because any parameter in $\Upsilon_1^* = \cup_{h=1}^{M_0} \Upsilon_{1h}^*$ can generate the true density $f_{M_0}(\mathbf{w}; \boldsymbol{\vartheta}_{M_0}^*) = \sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \gamma^*, \boldsymbol{\theta}_j^*)$, we need to restrict the estimators under the $(M_0 + 1)$ -component model to be in a neighborhood of Υ_{1h}^* in order to test $H_{0,1h}$.

Recall that $\mu_1^* < \mu_2^* \dots < \mu_{M_0}^*$. Let $\underline{\Theta}_\mu$ and $\overline{\Theta}_\mu$ denote the lower bound and upper bounds of Θ_μ . Define $D_1^* = [\underline{\Theta}_\mu, \frac{\mu_1^* + \mu_2^*}{2}] \times \Theta_\beta \times \Theta_{\sigma^2}$, $D_h^* = [\frac{\mu_{h-1}^* + \mu_h^*}{2}, \frac{\mu_h^* + \mu_{h+1}^*}{2}] \times \Theta_\beta \times \Theta_{\sigma^2}$ for $h = 2, \dots, M_0 - 1$, $D_{M_0}^* = [\frac{\mu_{M_0-1}^* + \mu_{M_0}^*}{2}, \overline{\Theta}_\mu] \times \Theta_\beta \times \Theta_{\sigma^2}$. Then $D_h^* \subset \Theta_\theta$ is a neighborhood containing θ_h^* but not θ_j^* for $j \neq h$. For $h = 1, \dots, M_0$, give a small positive constant $\epsilon > 0$, define a restricted parameter space $\Psi_h^* \subset \Theta_{\vartheta_{M_0+1}}(\epsilon)$ as

$$\Psi_h^* = \left\{ \begin{aligned} & \alpha_1, \dots, \alpha_{M_0+1} \in [\epsilon, 1 - \epsilon]; \sum_{j=1}^{M_0+1} \alpha_j = 1; \gamma \in \Theta_\gamma; \boldsymbol{\theta} \in \Theta_\theta : \boldsymbol{\theta}_j \in D_j^* \text{ for } j = 1, \dots, h-1; \\ & \boldsymbol{\theta}_h, \boldsymbol{\theta}_{h+1} \in D_h^*; \boldsymbol{\theta}_j \in D_{j-1}^* \text{ for } j = h+2, \dots, M_0+1. \end{aligned} \right\} \quad (20)$$

Note that $\Psi_h^* \cap \Upsilon_{1h}^* \neq \emptyset$ and $\Psi_h^* \cap \Upsilon_{1l}^* = \emptyset$ if $h \neq l$ while $\cup_{h=1}^{M_0} \Psi_h^* = \Theta_{\vartheta_{M_0+1}}(\epsilon)$.

Let $\hat{\Psi}_h^*$ and \hat{D}_h^* be consistent estimators of Ψ_h^* and D_h^* , which can be constructed from a consistent estimator of $\boldsymbol{\vartheta}_{M_0}^*$ in the M_0 -component model. We test $H_{0,1h} : \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1}$ by estimating the $(M_0 + 1)$ -component model under the restriction that $\boldsymbol{\vartheta}_{M_0+1} \in \hat{\Psi}_h^*$.

For $h = 1, 2, \dots, M_0$, define the ‘‘local’’ PMLE that maximizes the log-likelihood function of the $(M_0 + 1)$ -component model under the constrain that $\boldsymbol{\vartheta}_{M_0+1} \in \hat{\Psi}_h^*$ in (20) by

$$\hat{\boldsymbol{\vartheta}}_{M_0+1}^h = \arg \max_{\boldsymbol{\vartheta}_{M_0+1} \in \hat{\Psi}_h^*} L_{M_0+1,n}(\boldsymbol{\vartheta}_{M_0+1}) + \tilde{p}_n(\boldsymbol{\vartheta}_{M_0+1}),$$

where

$$L_{M,n}(\boldsymbol{\vartheta}_M) := \sum_{i=1}^n \log f_M(\mathbf{W}_i; \boldsymbol{\vartheta}_M) \quad \text{and} \quad \tilde{p}_n(\boldsymbol{\vartheta}_M) := \sum_{j=1}^M p_n(\sigma_j^2; \hat{\sigma}_{0,j}^2).$$

with

$$p_n(\sigma_j^2; \hat{\sigma}_{0,j}^2) := -a_n \{ \hat{\sigma}_{0,j}^2 / \sigma_j^2 + \log(\sigma_j^2 / \hat{\sigma}_{0,j}^2) - 1 \}, \quad (21)$$

where $\hat{\sigma}_{0,j}^2$ is a root- n consistent estimator of $\sigma_{0,j}^2$ from M_0 -component model under the null hypothesis. Because $\hat{\sigma}_j^2 - \sigma_{0,j}^2 = O_p(n^{-1/4})$ under the null hypothesis (c.f., Proposition 4(a)), $p_n(\hat{\sigma}_j^2; \hat{\sigma}_{0,j}^2) = o_p(1)$ when a_n is chosen to be $o(n^{1/4})$.

Under $H_0 : M = M_0$, Ψ_h^* contains a set of parameters Υ_{1h}^* defined in (19) such that $f_{M_0+1}(\mathbf{w}; \boldsymbol{\vartheta}_{M_0+1})$ is equal to $f_{M_0}(\mathbf{w}; \boldsymbol{\vartheta}_{M_0}^*)$ for any $\boldsymbol{\vartheta}_{M_0+1} \in \Upsilon_{1h}^*$ and is therefore the density function from which the data is generated. These penalized likelihood estimators are consistent.

Proposition 6. *Suppose that Assumption 1 holds. Then, under the null hypothesis $H_0 : M = M_0$, $\inf_{\boldsymbol{\vartheta}_{M_0+1} \in \Psi_h^*} |\hat{\boldsymbol{\vartheta}}_{M_0+1}^h - \boldsymbol{\vartheta}_{M_0+1}| \xrightarrow{p} 0$ for $h = 1, 2, \dots, M_0$.*

Consider the local PLRTS for testing $H_{0,1h} : \boldsymbol{\vartheta}_h = \boldsymbol{\vartheta}_{h+1}$ defined by

$$PLR_n^{M_0,h} := 2\{L_{M_0+1,n}(\hat{\boldsymbol{\vartheta}}_{M_0+1}^h) + \tilde{p}_n(\hat{\boldsymbol{\vartheta}}_{M_0+1}^h) - L_{M_0,n}(\hat{\boldsymbol{\vartheta}}_{M_0})\} \quad \text{for } h = 1, 2, \dots, M_0.$$

The test utilizing the local PLRTS, denoted as $PLR_n^{M_0,h}$, possesses power solely against local alternatives within the restricted parameter space of Ψ_h^* . To guarantee power against local alternatives over a wide range of directions, we consider the PLRTS characterized by the maximum of the local PLRTS for $h = 1, \dots, M_0$, as defined by

$$PLR_n(M_0) := \max\{PLR_n^{M_0,1}, PLR_n^{M_0,2}, \dots, PLR_n^{M_0,M_0}\}. \quad (22)$$

Because $\Theta_{\boldsymbol{\vartheta}_{M_0+1}}(\epsilon) = \cup_{h=1}^{M_0} \hat{\Psi}_h^*$, $PLR_n(M_0)$ is identical to $\max_{\boldsymbol{\vartheta}_{M_0+1} \in \Theta_{\boldsymbol{\vartheta}_{M_0+1}}(\epsilon)} \{L_{M_0+1,n}(\boldsymbol{\vartheta}_{M_0+1}) + \tilde{p}_n(\boldsymbol{\vartheta}_{M_0+1})\} - L_{M_0,n}(\hat{\boldsymbol{\vartheta}}_{M_0})$.

To derive the asymptotic null distribution of $PLR_n(M_0)$, collect the score vector for testing $H_{0,1h}$ for $h = 1, \dots, M_0$ into one vector as

$$\tilde{\mathbf{s}}(\mathbf{W}) = \begin{pmatrix} \tilde{\mathbf{s}}_\eta(\mathbf{W}) \\ \tilde{\mathbf{s}}_{\lambda\lambda}(\mathbf{W}) \end{pmatrix}, \quad \text{where } \tilde{\mathbf{s}}_\eta(\mathbf{W}) = \begin{pmatrix} \mathbf{s}_\alpha(\mathbf{W}) \\ \mathbf{s}_{(\gamma,\nu)}(\mathbf{W}) \end{pmatrix}_{(M_0+p+q+1) \times 1} \quad \text{and } \tilde{\mathbf{s}}_{\lambda\lambda}(\mathbf{W}) = \begin{pmatrix} s_{\lambda\lambda}^1(\mathbf{W}) \\ \vdots \\ s_{\lambda\lambda}^{M_0}(\mathbf{W}) \end{pmatrix}, \quad (23)$$

where

$$\begin{aligned} \mathbf{s}_\alpha(\mathbf{W}) &= \begin{pmatrix} f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_1^*) - f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_{M_0}^*) \\ \vdots \\ f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_{M_0-1}^*) - f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_{M_0}^*) \end{pmatrix} / f_{M_0}(\mathbf{W}; \boldsymbol{\vartheta}_{M_0}^*), \\ \mathbf{s}_{(\gamma,\nu)}(\mathbf{W}) &= \sum_{j=1}^{M_0} \alpha_j^* \nabla_{(\gamma,\nu)} f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_j^*) / f_{M_0}(\mathbf{W}; \boldsymbol{\vartheta}_{M_0}^*), \\ s_{\lambda\lambda}^h(\mathbf{W}) &= \tilde{\nabla}_{\boldsymbol{\theta}_h, \boldsymbol{\theta}_h^\top} f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_h^*) / f_{M_0}(\mathbf{W}; \boldsymbol{\vartheta}_{M_0}^*) \quad \text{for } h = 1, 2, \dots, M_0, \end{aligned} \quad (24)$$

with $\tilde{\nabla}_{\theta_h \theta_h^\top} f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}_h^*) := (c_{11} \nabla_{\theta_{h1} \theta_{h1}} f^*, \dots, c_{(q+2)(q+2)} \nabla_{\theta_{h,q+2} \theta_{h,q+2}} f^*, c_{12} \nabla_{\theta_{h1} \theta_{h2}} f^*, \dots, c_{(q+1)(q+2)} \nabla_{\theta_{h,q+1} \theta_{h,q+2}} f^*)^\top$ for $\boldsymbol{\theta}_h := (\theta_{h1}, \theta_{h2}, \theta_{h3}, \dots, \theta_{h,q+2})^\top := (\mu_h, \sigma_h^2, \beta_{h1}, \dots, \beta_{hq})^\top$ and $c_{jk} = 1/2$ for $j \neq k$ and $c_{jk} = 1$ for $j = k$. Define

$$\begin{aligned} \tilde{\mathcal{I}} &:= \mathbb{E}[\tilde{s}(\mathbf{W})\tilde{s}(\mathbf{W})^\top], \quad \tilde{\mathcal{I}}_\eta := \mathbb{E}[\tilde{s}_\eta(\mathbf{W})\tilde{s}_\eta(\mathbf{W})^\top], \quad \tilde{\mathcal{I}}_{\lambda\eta} := \mathbb{E}[\tilde{s}_{\lambda\lambda}(\mathbf{W})\tilde{s}_\eta(\mathbf{W})^\top], \\ \tilde{\mathcal{I}}_{\eta\lambda} &:= \tilde{\mathcal{I}}_{\lambda\eta}^\top, \quad \tilde{\mathcal{I}}_{\lambda\lambda} := \mathbb{E}[\tilde{s}_{\lambda\lambda}(\mathbf{W})\tilde{s}_{\lambda\lambda}(\mathbf{W})^\top], \quad \tilde{\mathcal{I}}_{\lambda,\eta} := \tilde{\mathcal{I}}_{\lambda\lambda} - \tilde{\mathcal{I}}_{\lambda\eta}\tilde{\mathcal{I}}_\eta^{-1}\tilde{\mathcal{I}}_{\eta\lambda}. \end{aligned} \quad (25)$$

Then, the asymptotic distribution of the normalized score function is given by

$$\tilde{\mathbf{S}}_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{s}(\mathbf{W}_i) \xrightarrow{d} \tilde{\mathbf{S}} \sim N(\mathbf{0}, \tilde{\mathcal{I}}),$$

where, in view of (23), $\tilde{\mathbf{S}}$ may be partitioned as $\tilde{\mathbf{S}} = (\tilde{\mathbf{S}}_\eta^\top, \tilde{\mathbf{S}}_{\lambda\lambda}^\top)^\top$ with $n^{-1/2} \sum_{i=1}^n \tilde{s}_\eta(\mathbf{W}_i) \xrightarrow{d} \tilde{\mathbf{S}}_\eta$ and $n^{-1/2} \sum_{i=1}^n \tilde{s}_{\lambda\lambda}(\mathbf{W}_i) \xrightarrow{d} \tilde{\mathbf{S}}_{\lambda\lambda}$.

Let $\tilde{\mathbf{S}}_{\lambda,\eta} := (\mathbf{S}_{\lambda,\eta}^1, \dots, \mathbf{S}_{\lambda,\eta}^{M_0})^\top := \tilde{\mathbf{S}}_{\lambda\lambda} - \tilde{\mathcal{I}}_{\lambda\eta}\tilde{\mathcal{I}}_\eta^{-1}\tilde{\mathbf{S}}_\eta \sim N(0, \tilde{\mathcal{I}}_{\lambda,\eta})$ be a $\mathbb{R}^{M_0(q+2)(q+1)/2}$ -valued random vector. For $h = 1, 2, \dots, M_0$, define $\tilde{\mathcal{I}}_{\lambda,\eta}^h := \mathbb{E}[\mathbf{S}_{\lambda,\eta}^h(\mathbf{S}_{\lambda,\eta}^h)^\top]$ and $\mathbf{G}_{\lambda,\eta}^h := (\tilde{\mathcal{I}}_{\lambda,\eta}^h)^{-1}\mathbf{S}_{\lambda,\eta}^h$.

Define $\hat{\mathbf{t}}_\lambda^h$ analogously to $\hat{\mathbf{t}}_\lambda$ as:

$$r_\lambda^h(\hat{\mathbf{t}}_\lambda^h) = \inf_{\mathbf{t}_\lambda^h \in \Lambda_\lambda} r^h(\mathbf{t}_\lambda^h); \quad r_\lambda^h(\mathbf{t}_\lambda^h) := (\mathbf{t}_\lambda^h - \mathbf{G}_{\lambda,\eta}^h)^\top \tilde{\mathcal{I}}_{\lambda,\eta}^h (\mathbf{t}_\lambda^h - \mathbf{G}_{\lambda,\eta}^h) \quad \text{for } h = 1, 2, \dots, M_0. \quad (26)$$

The local quadratic-form approximation of the log-likelihood function $LR_n^{M_0,h}$ around $\Upsilon_{1h}^* \subset \Theta_{\theta_{M_0+1}}$ shares an identical structure to the approximation we derived in Section 3 in testing H_{01} in the test of homogeneity. Consequently, we can show that $PLR_n^{M_0,h} \xrightarrow{d} (\hat{\mathbf{t}}_\lambda^h)^\top \tilde{\mathcal{I}}_{\lambda,\eta}^h \hat{\mathbf{t}}_\lambda^h$. Then, given (22), the asymptotic null distribution of the PLRTS for testing H_{01} is given by the maximum over $(\hat{\mathbf{t}}_\lambda^h)^\top \tilde{\mathcal{I}}_{\lambda,\eta}^h \hat{\mathbf{t}}_\lambda^h$'s for $h = 1, 2, \dots, M_0$.

Assumption 3. (a) $\alpha_j^* \in (\epsilon, 1 - \epsilon)$ for $j = 1, \dots, M_0$. (b) $\tilde{\mathcal{I}}$ is non-singular. (c) a_n in (21) satisfies $a_n = O(1)$.

Proposition 7. Suppose that Assumptions 1-3 are satisfied. Then under the null hypothesis $H_0 : M = M_0$, $PLR_n(M_0) \xrightarrow{d} \max\{(\hat{\mathbf{t}}_\lambda^1)^\top \tilde{\mathcal{I}}_{\lambda,\eta}^1 \hat{\mathbf{t}}_\lambda^1, \dots, (\hat{\mathbf{t}}_\lambda^{M_0})^\top \tilde{\mathcal{I}}_{\lambda,\eta}^{M_0} \hat{\mathbf{t}}_\lambda^{M_0}\}$.

The asymptotic null distribution of $PLR_n(M_0)$ is non-standard but it is straightforward to simulate the random variable from the asymptotic null distribution using the estimates. Specifically, we simulate a draw of $\tilde{\mathbf{S}}_{\lambda,\eta} = (\mathbf{S}_{\lambda,\eta}^1, \dots, \mathbf{S}_{\lambda,\eta}^{M_0})^\top$ from $N(0, \hat{\tilde{\mathcal{I}}}_{\lambda,\eta})$, where $\hat{\tilde{\mathcal{I}}}_{\lambda,\eta}$ is a sample analogue estimator of $\tilde{\mathcal{I}}_{\lambda,\eta}$. Then, compute $\mathbf{G}_{\lambda,\eta}^h = (\hat{\tilde{\mathcal{I}}}_{\lambda,\eta}^h)^{-1}\mathbf{S}_{\lambda,\eta}^h$ and obtain $\hat{\mathbf{t}}_\lambda^h$ analogously to (26) using an estimator of $\tilde{\mathcal{I}}_{\lambda,\eta}^h$ for $h = 1, \dots, M_0$, and a simulated random draw is computed as $\max\{(\hat{\mathbf{t}}_\lambda^1)^\top \hat{\tilde{\mathcal{I}}}_{\lambda,\eta}^1 \hat{\mathbf{t}}_\lambda^1, \dots, (\hat{\mathbf{t}}_\lambda^{M_0})^\top \hat{\tilde{\mathcal{I}}}_{\lambda,\eta}^{M_0} \hat{\mathbf{t}}_\lambda^{M_0}\}$. Appendix B.2-B.3 present an expression for score functions using Hermit polynomials.

5 EM test for $H_0 : M = M_0$ against $H_A : M = M_0 + 1$

This section develops an EM test used for testing the hypothesis $H_0 : M = M_0$ against the alternative hypothesis $H_A : M = M_0 + 1$. A key limitation of the PLRT, as discussed in the previous section, is that the computation of mixing probabilities, denoted as α_j , is subject to a hard constraint, which is dictated by an arbitrary choice of bounds. The EM test, on the other hand, circumvents the need for imposing an explicit constraint on the α_j values. It achieves this by performing a limited number of EM steps, starting from a predetermined set of α_j values. The EM test approach offers certain advantages, including computational simplicity and less stringent assumptions.

Let \mathcal{T} be a finite set of numbers in $(0, 0.5]$ with $0.5 \in \mathcal{T}$, and let $p(\tau) \leq 0$ be a penalty term that is continuous in τ , $p(0.5) = 0$, and $p(\tau) \rightarrow -\infty$ as τ goes to 0. Specifically, we choose

$$p(\tau) := \log(2 \min\{\tau, 1 - \tau\}).$$

For each $\tau_0 \in \mathcal{T}$, let $\tau^{(1)}(\tau_0) = \tau_0$ and define the restricted penalized MLE by

$$\boldsymbol{\vartheta}_{M_0+1}^{h(1)}(\tau_0) = \arg \max_{\boldsymbol{\vartheta}_{M_0+1} \in \Theta_{\boldsymbol{\vartheta}_{M_0+1}}^h(\tau)} PL_n(\boldsymbol{\vartheta}_{M_0+1}, \tau_0)$$

where $\Theta_{\boldsymbol{\vartheta}_{M_0+1}}^h(\tau_0) := \{\boldsymbol{\theta} \in \hat{\Psi}_h : \alpha_h / (\alpha_h + \alpha_{h+1}) = \tau_0\}$ and

$$PL_n(\boldsymbol{\vartheta}_{M_0+1}, \tau) := L_{M_0+1,n}(\boldsymbol{\vartheta}_{M_0+1}) + \tilde{p}_n(\boldsymbol{\vartheta}_{M_0+1}) + p(\tau).$$

Starting from $(\boldsymbol{\vartheta}_{M_0+1}^{h(1)}(\tau_0), \tau^{h(1)}(\tau_0))$ with $\tau^{h(1)}(\tau_0) = \tau_0$, update $\boldsymbol{\vartheta}_{M_0+1}^{h(k)}(\tau_0)$ and $\tau^{h(k)}(\tau_0)$ by the following generalized EM algorithm. Denote the estimators after the k -th round of EM algorithm iteration by $\boldsymbol{\vartheta}_{M_0+1}^{h(k)}$ and $\tau^{h(k)}$. In the E-step, for $i = 1, \dots, N$ and $j = 1, \dots, M_0 + 1$, compute the weight for observation i and type j as:

$$\begin{aligned} w_{ij}^{(k)} &= \begin{cases} \alpha_j^{(k)} f(\mathbf{W}_i; \boldsymbol{\gamma}^{(k)}, \boldsymbol{\theta}_j^{(k)}) / f_{M_0+1}(\mathbf{W}_i; \boldsymbol{\vartheta}_{M_0+1}^{h(k)}(\tau_0)), & j = 1, \dots, h-1, \\ \alpha_{j-1}^{(k)} f(\mathbf{W}_i; \boldsymbol{\gamma}^{(k)}, \boldsymbol{\theta}_j^{(k)}) / f_{M_0+1}(\mathbf{W}_i; \boldsymbol{\vartheta}_{M_0+1}^{h(k)}(\tau_0)), & j = h+2, \dots, M_0+1, \end{cases} \\ w_{ih}^{(k)} &= \tau^{h(k)} \alpha_h^{(k)} f(\mathbf{W}_i; \boldsymbol{\gamma}^{(k)}, \boldsymbol{\theta}_h^{(k)}) / f_{M_0+1}(\mathbf{W}_i; \boldsymbol{\vartheta}_{M_0+1}^{h(k)}(\tau_0)), \\ w_{i,h+1}^{(k)} &= (1 - \tau^{h(k)}) \alpha_h^{(k)} f(\mathbf{W}_i; \boldsymbol{\gamma}^{(k)}, \boldsymbol{\theta}_{h+1}^{(k)}) / f_{M_0+1}(\mathbf{W}_i; \boldsymbol{\vartheta}_{M_0+1}^{h(k)}(\tau_0)), \end{aligned} \tag{27}$$

where, for brevity, we drop the superscript h and its dependency on τ_0 from the notations such as $w_{ij}^{h(k)}(\tau_0)$.

In the M-step, we update α and τ by

$$\alpha_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n w_{ij}^{(k)} \quad \text{for } j = 1, \dots, M_0 + 1 \quad \text{and}$$

$$\tau^{h(k+1)} = \arg \min_{\tau} \left\{ \sum_{i=1}^n w_{ih}^{(k)} \log(\tau) + \sum_{i=1}^n w_{i,h+1}^{(k)} \log(1 - \tau) + p(\tau) \right\}.$$

We also update θ_j and γ as

$$(\sigma_j^{(k+1)})^2 = \arg \min_{\sigma_j^2} \left\{ \sum_{i=1}^n w_i^{j(k)} \sum_{t=1}^T (y_{it} - \mu_j^{(k+1)} - \mathbf{z}_{it}^\top \boldsymbol{\gamma}^{(k+1)} - \mathbf{x}_{it}^\top \boldsymbol{\beta}_j^{(k+1)})^2 + p_n(\sigma_j^2) \right\},$$

$$\boldsymbol{\gamma}^{(k+1)} = \left(\sum_{i=1}^n \sum_{t=1}^T \mathbf{z}_{it} \mathbf{z}_{it}^\top \right)^{-1} \left(\sum_{i=1}^n \sum_{t=1}^T \mathbf{z}_{it} \left(y_{it} - \sum_{j=1}^{M_0+1} w_{ij}^{(k)} \tilde{\mathbf{x}}_{it}^\top \begin{pmatrix} \mu_j^{(k)} \\ \boldsymbol{\beta}_j^{(k)} \end{pmatrix} \right) \right), \quad \text{and}$$

$$\begin{pmatrix} \mu_j^{(k+1)} \\ \boldsymbol{\beta}_j^{(k+1)} \end{pmatrix} = \left(\sum_{i=1}^n w_{ij}^{(k)} \sum_{t=1}^T \tilde{\mathbf{x}}_{it} \tilde{\mathbf{x}}_{it}^\top \right)^{-1} \left(\sum_{i=1}^n w_{ij}^{(k)} \sum_{t=1}^T \tilde{\mathbf{x}}_{it} (y_{it} - \mathbf{z}_{it}^\top \boldsymbol{\gamma}^{(k+1)}) \right),$$

where $\tilde{\mathbf{x}}_{it} = (1, \mathbf{x}_{it}^\top)^\top$. In the updating procedure, $\boldsymbol{\vartheta}_{M_0+1}^{h(k+1)}(\tau_0)$ is not restricted to be in $\hat{\boldsymbol{\Psi}}_h^*$.

For each $\tau_0 \in \mathcal{T}$ and each step k , define

$$M_n^{h(k)}(\tau_0) := 2 \left\{ PL_n(\boldsymbol{\vartheta}_{M_0+1}^{h(k)}(\tau_0), \tau^{h(k)}(\tau_0)) - L_{M_0, n}(\hat{\boldsymbol{\vartheta}}_{M_0}) \right\}. \quad (28)$$

With a pre-determined finite number K , define the *local* EM test statistic by taking maximum of $M_n^{h(k)}(\tau_0)$ across different τ_0 's as

$$EM_n^h := \max\{M_n^{h(K)}(\tau_0) : \tau_0 \in \mathcal{T}\}. \quad (29)$$

The test statistic EM_n^h tests $H_{0,1h} : \boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1}$ and has a power against the local alternative that splits the h -th component of the null M_0 -component model into two different components. To achieve power against a wide range of local alternatives, we consider the EM test statistic that takes the maximum of M_0 local EM test statistics:

$$EM_n(M_0) := \max\{EM_n^{1(K)}, \dots, EM_n^{M_0(K)}\}. \quad (30)$$

Proposition 8. *Suppose that Assumptions 1–3 hold and $\{0.5\} \in \mathcal{T}$. Then, under the null hypothesis $H_0 : M = M_0$, for any finite K , $EM_n(M_0) \xrightarrow{d} \max\{(\hat{\mathbf{t}}_\lambda^1)^\top \boldsymbol{\mathcal{I}}_{\lambda, \eta}^1 \hat{\mathbf{t}}_\lambda^1, \dots, (\hat{\mathbf{t}}_\lambda^{M_0})^\top \boldsymbol{\mathcal{I}}_{\lambda, \eta}^{M_0} \hat{\mathbf{t}}_\lambda^{M_0}\}$.*

Therefore, the asymptotic null distribution of EM test statistic $EM_n(M_0)$ is the same as that of the PLRTS.

6 Asymptotic Distribution under Local Alternatives

We derive the asymptotic distribution of PLRTS and EM test static under local alternatives. For brevity, we focus on testing $H_0 : M = 1$ against $H_A : M = 2$. Consider the following local alternative to the homogeneous model $f(\mathbf{w}; \gamma^*, \boldsymbol{\theta}^*)$ with $\boldsymbol{\theta}^* = (\mu^*, \sigma^{*2}, (\boldsymbol{\beta}^*)^\top)^\top$. For brevity, we omit the common parameter γ in this section. In a reparameterized parameter, $\boldsymbol{\psi}^* = ((\boldsymbol{\nu}^*)^\top, (\boldsymbol{\lambda}^*)^\top)^\top$. For $\alpha^* \in (0, 1)$ and a local parameter $\mathbf{h} = (\mathbf{h}_\nu^\top, \mathbf{h}_\lambda^\top)^\top$ with $\mathbf{h}_\lambda \in v(\Theta_\lambda)$, we consider a sequence of contiguous local alternatives $(\alpha_n, \boldsymbol{\psi}_n^\top)^\top = (\alpha_n, \boldsymbol{\nu}_n^\top, \boldsymbol{\lambda}_n^\top) \in \Theta_\alpha \times \Theta_\nu \times \Theta_\lambda$ such that, with $\mathbf{t}_\lambda(\boldsymbol{\lambda}, \alpha)$ given by (10),

$$\mathbf{h}_\nu = \sqrt{n}(\boldsymbol{\nu}_n - \boldsymbol{\nu}^*), \quad \mathbf{h}_\lambda = \sqrt{n}\mathbf{t}_\lambda(\boldsymbol{\lambda}_n, \alpha_n), \quad \text{and} \quad \alpha_n = \alpha^* + o(1). \quad (31)$$

Equivalently, the non-reparameterized contiguous local alternatives are given by

$$\boldsymbol{\theta}_{1,n} = \boldsymbol{\nu}_n + (1 - \alpha_n)\boldsymbol{\lambda}_n \quad \text{and} \quad \boldsymbol{\theta}_{2,n} = \boldsymbol{\nu}_n - \alpha_n\boldsymbol{\lambda}_n \quad (32)$$

for $\boldsymbol{\nu}_n = \boldsymbol{\nu}^* + n^{-1/2}\mathbf{h}_\nu$ and $\boldsymbol{\lambda}_n = (\lambda_{1,n}, \lambda_{2,n}, \dots, \lambda_{q+2,n})^\top$ with

$$\lambda_{j,n} = n^{-1/4}(\alpha_n(1 - \alpha_n))^{-1/2}h_{\lambda,j} \quad \text{for } j = 1, \dots, q + 2,$$

where $\mathbf{h}_\lambda = (h_{\lambda,1}^2, \dots, h_{\lambda,q+2}^2, h_{\lambda,1}h_{\lambda,2}, \dots, h_{\lambda,q+1}h_{\lambda,q+2})^\top$. The local alternatives are of order $n^{1/4}$ rather than $n^{1/2}$. See the discussion after Proposition 4.

The following proposition provides the asymptotic distribution of the PLRT and EM test statistics under contiguous local alternatives.

Proposition 9. *Suppose that the assumptions in Proposition 8 hold for $M_0 = 1$. Consider a sequence of contiguous local alternatives $\boldsymbol{\vartheta}_{2,n} = (\alpha_n, \boldsymbol{\theta}_{1,n}^\top, \boldsymbol{\theta}_{2,n}^\top)^\top$ given in (32), where α_n and $\boldsymbol{\lambda}_n$ satisfy (31). Then, under $H_{1,n} : \boldsymbol{\vartheta} = \boldsymbol{\vartheta}_{2,n}$, we have $PLR_n(1), EM_n(1) \xrightarrow{d} (\tilde{\mathbf{t}}_\lambda)^\top \mathcal{I}_{\lambda,\eta} \tilde{\mathbf{t}}_\lambda$, where $\tilde{\mathbf{t}}_\lambda$ has the same distribution as $\hat{\mathbf{t}}_\lambda$ in Proposition 4 but replacing $\mathbf{G}_{\lambda,\eta}$ with $(\mathcal{I}_{\lambda,\eta})^{-1}\mathbf{S}_{\lambda,\eta} + \mathbf{h}_\lambda$.*

Importantly, a set of contiguous local alternatives considered in (32) excludes a sequence such that $\alpha_n \rightarrow 0$ or 1.

7 Sequential Hypothesis Testing

To estimate the number of components, we sequentially test $H_0 : M = r$ against $H_1 : M = r + 1$ starting from $r = 1$, and then $r = 2, \dots, \bar{M}$, where \bar{M} is the upper bound for the number of components, which is assumed to be larger than M_0 . The first value for r that leads to a nonrejection

of H_0 gives our estimate for M_0 . Robin and Smith (2000) develops a similar sequential hypothesis test for estimating the rank of a matrix.

For $M = 1, \dots, \bar{M}$, let $c_{1-q_n}^M$ denote the $100(1 - q_n)$ percentile of the cumulative distribution function of a random variable $\max\{(\hat{\mathbf{t}}_\lambda^1)^\top \mathcal{I}_{\lambda, \eta}^1 \hat{\mathbf{t}}_\lambda^1, \dots, (\hat{\mathbf{t}}_\lambda^M)^\top \mathcal{I}_{\lambda, \eta}^M \hat{\mathbf{t}}_\lambda^M\}$ for $M = M_0$ in Propositions 7 and 8. Let $\hat{c}_{1-q_n}^M$ be a consistent estimator of $c_{1-q_n}^M$. Then, our estimator based on sequential hypothesis testing (SHT, hereafter) is defined as

$$\begin{aligned}\hat{M}_{\text{PLR}} &= \min_{M \in \{0, \dots, \bar{M}\}} \{M : \text{PLR}_n(r) \geq \hat{c}_{1-q_n}^r, r = 0, \dots, M-1, \text{PLR}_n(M) < \hat{c}_{1-q_n}^M\}, \\ \hat{M}_{\text{EM}} &= \min_{M \in \{0, \dots, \bar{M}\}} \{M : \text{EM}_n(r) \geq \hat{c}_{1-q_n}^r, r = 0, \dots, M-1, \text{EM}_n(M) < \hat{c}_{1-q_n}^M\}.\end{aligned}\quad (33)$$

The estimators \hat{M}_{PLR} and \hat{M}_{EM} depend on the choice of the significance level q_n . The following proposition states that \hat{M}_{PLR} and \hat{M}_{EM} converge to M_0 in probability as $n \rightarrow \infty$ if we choose q_n such that q_n such that $-n^{-1} \ln q_n = o(1)$ and $q_n = o(1)$.

Let $Q_n^M(\boldsymbol{\vartheta}_M) := n^{-1} \sum_{i=1}^n \ln f_M(\mathbf{w}_i; \boldsymbol{\vartheta}_M)$ and $Q^M(\boldsymbol{\vartheta}_M) := \mathbb{E}[\ln f_M(\mathbf{w}_i; \boldsymbol{\vartheta}_M)]$, where $f_M(\mathbf{w}_i; \boldsymbol{\vartheta}_M)$ is defined in (1) for $M = 1, \dots, \bar{M}$.

Assumption 4. For $M = 1, \dots, M_0 - 1$, (a) $Q^M(\boldsymbol{\vartheta}_M)$ has a unique maximum at $\boldsymbol{\vartheta}_M^*$ in $\Theta_{\boldsymbol{\vartheta}_M}$; (b) $\Theta_{\boldsymbol{\vartheta}_M}$ is compact; (c) $\boldsymbol{\vartheta}_M^*$ is interior to $\Theta_{\boldsymbol{\vartheta}_M}$; (d) $B^M(\boldsymbol{\vartheta}_M^*) := \mathbb{E} \left\{ \nabla_{\boldsymbol{\vartheta}_M} \ln f_M(\mathbf{w}_i; \boldsymbol{\vartheta}_M) \nabla_{\boldsymbol{\vartheta}_M}^\top \ln f_M(\mathbf{w}_i; \boldsymbol{\vartheta}_M) \right\}$ is nonsingular; (e) $A^M(\boldsymbol{\vartheta}_M^*) := \mathbb{E} \left\{ \nabla_{\boldsymbol{\vartheta}_M} \ln f_M(\mathbf{w}_i; \boldsymbol{\vartheta}_M) \right\}$ has constant rank in some open neighborhood of $\boldsymbol{\vartheta}_M^*$; (f) $Q^{M+1}(\boldsymbol{\vartheta}_{M+1}^*) - Q^M(\boldsymbol{\vartheta}_M^*) > 0$.

Proposition 10. Suppose that $M_0 < \bar{M}$ and Assumptions 1-4 hold. If we choose q_n such that $-n^{-1} \ln q_n = o(1)$ and $q_n = o(1)$, then $\hat{M}_{\text{PLR}} - M_0 = o_p(1)$ and $\hat{M}_{\text{EM}} - M_0 = o_p(1)$.

Assumption 4 (a)-(e) ensure the consistency and asymptotic normality of $\hat{\boldsymbol{\vartheta}}_M$, where (c)-(e) correspond to Assumption A6 of White (1982). Per Assumption 4(f), the Kullback-Leibler Information Criterion of the model relative to the true M_0 components model strictly decreases as the number of components M increases, for $M < M_0$.

8 Simulation

In this section, we examine the finite sample performance of the EM test and PLRT by simulation. We test $H_0 : M = M_0$ against $H_1 : M = M_0 + 1$ for the model with $M_0 = 2$ and 3.

8.1 Choice of penalty function

We have developed a data-dependent empirical formula for a_n by selecting a formula that ensures empirical rejection probabilities match the nominal size (5%) across various null models and sam-

ple sizes, as reported in Table 11 in Appendix D. Specifically, for the model without conditioning variables, we have derived the following data-dependent empirical formula for testing the null hypotheses of $M_0 = 1, 2, 3, 4$:

$$a_n = \begin{cases} \left(1 + \exp \left\{ \frac{\hat{\rho}_1^{M_0}}{\hat{\rho}_4^{M_0}} + \frac{\hat{\rho}_2^{M_0}}{\hat{\rho}_4^{M_0}} \frac{1}{T} + \frac{\hat{\rho}_3^{M_0}}{\hat{\rho}_4^{M_0}} \frac{1}{n} \right\} \right)^{-1}, & M_0 = 1 \\ \left(1 + \exp \left\{ \frac{\hat{\rho}_1^{M_0}}{\hat{\rho}_4^{M_0}} + \frac{\hat{\rho}_2^{M_0}}{\hat{\rho}_4^{M_0}} \frac{1}{T} + \frac{\hat{\rho}_3^{M_0}}{\hat{\rho}_4^{M_0}} \frac{1}{n} + \frac{\hat{\rho}_5^{M_0}}{\hat{\rho}_4^{M_0}} \log \left(\frac{\omega(\boldsymbol{\vartheta}_{M_0}; M_0)}{1 - \omega(\boldsymbol{\vartheta}_{M_0}; M_0)} \right) \right\} \right)^{-1}, & M_0 = 2, 3, 4, \end{cases} \quad (34)$$

where $\omega(\boldsymbol{\vartheta}_{M_0}; M_0)$ is the misclassification probability as defined in Melnykov and Maitra (2010) for each of the null models. The parameters $\hat{\rho}_1^{M_0}$, $\hat{\rho}_2^{M_0}$, $\hat{\rho}_3^{M_0}$, $\hat{\rho}_4^{M_0}$, and $\hat{\rho}_5^{M_0}$ are chosen as follows. Across different null models, sample sizes, and various candidate values of a_n , we estimate the empirical rejection probabilities at the 5% significance level by simulations and denote them by \hat{s} . For example, when testing $H_0 : M_0 = 2$, we repeatedly simulate the 500 datasets under each of the 48 null model parameters and sample sizes $(N, T, \alpha, \mu, \sigma) \in \{100, 500\} \times \{2, 5, 10\} \times \{(0.5, 0.5), (0.2, 0.8)\} \times \{(-1, 1), (-0.5, 0.5), (-0.5, 0.8)\} \times \{(1, 1), (1.5, 0.75), (0.8, 1.2)\}$ and test the null hypothesis of $H_0 : M_0 = 2$ by the EM test using one of the six values of $a_n \in \{0.01, 0.05, 0.1, 0.2, 0.3, 0.4\}$. For each of $108 \times 6 = 648$ combinations of the parameter values, sample sizes, and a_n values, let \hat{s} denote the fraction of simulated datasets that led to the rejection of the null hypothesis at 5% significance level. Using these 648 ‘‘observations’’ of $\{\hat{s}, N, T, \omega(\boldsymbol{\vartheta}_2; 2), a_n\}$, we run the following regression:

$$\begin{aligned} & \log \left(\frac{\hat{s}}{1 - \hat{s}} \right) - \log \left(\frac{0.05}{1 - 0.05} \right) \\ &= \begin{cases} \rho_1^{M_0} + \rho_2^{M_0} \frac{1}{T} + \rho_3^{M_0} \frac{1}{n} + \rho_4^{M_0} \log \left(\frac{a_n}{1 - a_n} \right), & M_0 = 1 \\ \rho_1^{M_0} + \rho_2^{M_0} \frac{1}{T} + \rho_3^{M_0} \frac{1}{n} + \rho_4^{M_0} \log \left(\frac{a_n}{1 - a_n} \right) + \rho_5^{M_0} \log \left(\frac{\omega(\boldsymbol{\vartheta}_{M_0}; M_0)}{1 - \omega(\boldsymbol{\vartheta}_{M_0}; M_0)} \right), & M_0 = 2, 3, 4, \end{cases} \end{aligned}$$

where $\hat{\rho}_1^{M_0}$, $\hat{\rho}_2^{M_0}$, $\hat{\rho}_3^{M_0}$, $\hat{\rho}_4^{M_0}$, and $\hat{\rho}_5^{M_0}$ in (34) denotes the corresponding estimates. Table 12 in the Appendix reports the estimates. Note that the data-dependent formula (34) is obtained by setting $\hat{s} = 0.05$ and solving for a_n in the above equation.

For the model with conditioning variables, we find that the value of a_n that gives accurate Type I errors is sensitive to the dimension of covariates, and developing a data-dependent empirical formula for a_n is difficult. Consequently, we choose a constant value of a_n that depends only on the number of components $M_0 = 1, 2, 3$, and 4 as follows: $a_n = 0.1617$ if $M_0 = 1$; $a_n = 0.0025$ if $M_0 = 2$; $a_n = 0.0567$ if $M_0 = 3$; $a_n = 0.4858$ if $M_0 = 4$; $a_n = 0.5$ if $M_0 \geq 5$. These penalty terms for the regression with covariates are chosen by averaging the prediction of the penalty function for the null parameters used in the simulations. For example, the penalty term for $M_0 = 2$ is chosen by generating a_n using the formula for all the combinations of $(N, T, \alpha, \mu, \sigma)$ in Table 11 for $M_0 = 2$

and taking the average across the predicted \hat{a}_n 's. For $M_0 \geq 5$, we use the parametric bootstrap method to obtain the critical values for our empirical application, where we set $a_n = 0.5$.

8.2 Simulation results

Table 1 displays the simulated Type I error rates for the EM test when examining the null hypothesis $H_0 : M = 2$ against the alternative hypothesis $H_1 : M = 3$. A total of 2000 repetitions were employed for the asymptotic distribution, while 1000 repetitions were used for the bootstrap distribution. Moreover, the PLRT with simulated critical values was considered.

The table presents results for four distinct null models, as explained in the table's footnote. Utilizing the asymptotic distribution, the EM test sizes generally approximate the nominal 5% level. Nonetheless, the test may be undersized in instances where $T \geq 5$. Furthermore, the test size is larger when the mixing proportions are equal ($\alpha = (0.5, 0.5)$) in comparison to when they are unequal ($\alpha = (0.2, 0.8)$). The bootstrapped EM test demonstrates satisfactory performance.

For the PLRT, 2000 repetitions were conducted, and results were reported for cases where a constraint was applied to $\alpha_j \in [\epsilon, 1 - \epsilon]$ with $\epsilon = 0.1$. The value of a_n for the PLRT was chosen to be ten times larger than its value for the EM test. The findings suggest that the PLRT is slightly oversized.

Table 2 reports the rejection frequency of testing $H_0 : M_0 = 2$ under 12 alternative three-component mixture models, as elaborated in the table's footnote. For both EM test and PLRT, the test power is greater when distances between μ_j 's are larger and equal, such as $(\mu_1, \mu_2, \mu_3) = (-1, 0, 1)$ or $(-1.5, 0, 1.5)$, as opposed to unbalanced distances like $(-1, 0, 2)$ or $(-0.5, 0, 1.5)$. The power is also improved when the mixture probabilities are equal ($\alpha = (1/3, 1/3, 1/3)$) rather than unequal ($\alpha = (1/4, 1/2, 1/4)$). The power increases with both the time-dimension T and cross-sectional sample size N . Reflecting a larger actual rejection frequency of the PLRT under $H_0 : M_0 = 2$ in Table 1, the power of the PLRT is often higher than that of the EM test, although the EM test sometimes has higher power, especially when the mixing probabilities are unequal.

Table 3 displays the simulated Type I error rates of the EM test using the asymptotic distribution for testing $H_0 : M_0 = 3$ against $H_1 : M_0 = 4$. Six null models are considered with varying $(\alpha_1, \alpha_2, \alpha_3)$ and (μ_1, μ_2, μ_3) values. The EM test generally yields accurate Type I errors.

The Type I error rates of the EM test with conditioning variables under the null $M_0 = 2$ are examined using 500 repetitions. Results presented in Table 4 indicate a slightly oversized test for small samples with $(N, T) = (200, 2)$, but overall, the finite sample properties are satisfactory.

In our empirical application examining production function heterogeneity in Japan and Chile, we find evidence that the number of components is frequently greater than five when we sequentially apply our EM test to estimate the number of components. Consequently, we also investigate the performance of the sequential hypothesis testing (SHT) using the EM test in comparison to the

AIC and the BIC when the data is generated from a five-component model in a realistic setting. Specifically, we simulate 100 datasets from the estimated five-component model of the Chilean textile industry in our empirical application and apply these three methods to select the number of components in each of the 100 datasets. Here, we apply the EM test at the 5 percent significance level to sequentially test the null hypothesis $H_0 : M = M_0$ for $M_0 = 1, 2, \dots, 7$, and we determine the number of components to be M_0 when we fail to reject $H_0 : M = M_0$ as in (33).

Table 5 presents the frequency at which the three methods select the number of components in this simulation. The table demonstrates that the proposed sequential hypothesis test selects the correct number of components 72 % of the time, while it underestimates the true number of components 25 % of the time. Conversely, the AIC overestimates the number of components 86 % of the time, and the BIC underestimates the number of components by selecting a four-component model 41 % of the time, accurately estimating the number of components 58% of the time. Overall, in this simulation, our proposed sequential hypothesis testing approach outperforms both the AIC and the BIC.

Table 1: Sizes (in %) of EM test and PLRT of $H_0 : M_0 = 2$ against $H_A : M_0 = 3$ at the 5% level

T	EM Test						EM Test						PLRT					
	Asymptotic						Parametric Bootstrap						Asymptotic					
	3		5		8		3		5		8		3		5		8	
N	200	400	200	400	200	400	200	400	200	400	200	400	200	400	200	400	200	400
(A, C)	5.3	4.8	4	3.8	4	2.95	4.6	6.2	6.2	4.6	4.8	5.2	7.7	6.6	6.75	6.7	6.5	5.5
(A, D)	5.9	4.9	5	5	4.45	4	5.4	4.8	5.2	5.6	5	5.8	5.4	5.15	5.1	6.1	5.55	6.2
(B, C)	3.8	2.5	3.45	3.05	3.6	3.25	3.6	5.6	4.2	5.2	4	5.4	6.25	5.45	6.45	6.05	5.05	6
(B, D)	4.8	4.6	3.5	3.15	3.55	3.95	3.6	3.6	5.8	4	6.2	4.6	2.35	4.4	3.9	4.85	4.95	5.2

¹ A and B refer to respectively $(\alpha_1, \alpha_2) = (0.5, 0.5)$ and $(0.2, 0.8)$, while C and D refer to $(\mu_1, \mu_2) = (-1, 1)$ and $(-0.5, 0.5)$, respectively.

² The variance is set to $(\sigma_1, \sigma_2) = (0.8, 1.2)$. The asymptotic simulations are based on 2000 repetitions and the bootstrap simulation is based on 1000 repetitions.

Table 2: Powers (in %) of EM test and PLRT of $H_0 : M_0 = 2$ against $H_A : M_0 = 3$ at the 5% level

N	A				B			
	100		500		100		500	
T	2	5	2	5	2	5	2	5
	EM test							
(C, G)	20.9	81.6	57.6	100.0	20.5	82.7	62.6	100.0
(C, H)	49.2	99.9	99.9	100.0	38.4	98.7	98.8	100.0
(C, I)	12.1	20.4	18.0	62.6	10.6	20.4	16.8	65.8
(D, G)	77.9	100.0	100.0	100.0	86.5	100.0	100.0	100.0
(D, H)	57.4	100.0	100.0	100.0	42.8	100.0	100.0	100.0
(D, I)	16.0	59.5	31.8	99.9	13.8	70.8	40.3	100.0
(E, G)	93.0	100.0	100.0	100.0	94.0	100.0	100.0	100.0
(E, H)	83.8	100.0	100.0	100.0	70.7	100.0	100.0	100.0
(E, I)	25.7	97.0	80.2	100.0	30.7	96.8	83.1	100.0
(F, G)	99.9	100.0	100.0	100.0	100.0	100.0	100.0	100.0
(F, H)	93.5	100.0	100.0	100.0	85.3	100.0	100.0	100.0
(F, I)	40.8	99.9	98.2	100.0	52.1	100.0	99.5	100.0
	PLRT							
(C, G)	22.7	85.1	56.5	100.0	23.2	82.6	58.7	100.0
(C, H)	57.1	100.0	99.8	100.0	43.2	99.7	99.1	100.0
(C, I)	12.0	21.0	12.4	66.1	11.3	22.1	12.4	69.3
(D, G)	79.9	100.0	100.0	100.0	87.6	100.0	100.0	100.0
(D, H)	65.3	100.0	100.0	100.0	49.1	100.0	100.0	100.0
(D, I)	14.8	63.6	28.6	100.0	13.6	75.2	36.7	100.0
(E, G)	91.5	100.0	100.0	100.0	93.7	100.0	100.0	100.0
(E, H)	86.8	100.0	100.0	100.0	75.9	100.0	100.0	100.0
(E, I)	28.7	97.2	77.7	100.0	33.0	97.9	85.5	100.0
(F, G)	99.9	100.0	100.0	100.0	100.0	100.0	100.0	100.0
(F, H)	96.4	100.0	100.0	100.0	89.5	100.0	100.0	100.0
(F, I)	45.7	100.0	98.4	100.0	57.2	100.0	99.8	100.0

Notes: A and B refer to $(\alpha_1, \alpha_2, \alpha_3) = (1/3, 1/3, 1/3)$ and $(1/4, 1/2, 1/4)$, respectively; $C, D, E,$ and F refer to $(\mu_1, \mu_2, \mu_3) = (-0.5, 0, 1.5), (-1, 0, 1), (-1, 0, 2), (-1.5, 0, 1.5)$, respectively; G, H, I refer to $(\sigma_1, \sigma_2, \sigma_3) = (0.6, 0.6, 1.2), (0.6, 1.2, 0.6), (1, 1, 1)$.

Table 3: Sizes (in %) of **EM test** of $H_0 : M_0 = 3$ against $H_A : M_0 = 4$ at 5% level

	(A,C)	(A,D)	(A,E)	(B,C)	(B,D)	(B,E)
100,2	5.95	5.15	5.05	5.05	5.85	4.40
500,2	5.60	5.55	5.25	5.10	5.65	4.05
100,5	4.30	6.00	4.20	5.15	5.10	5.70
500,5	4.20	4.55	3.95	4.50	4.15	4.15

Notes: A and B refer to $(\alpha_1, \alpha_2, \alpha_3) = (1/3, 1/3, 1/3)$ and $(0.25, 0.5, 0.25)$, respectively, while C, D, E refer to $(\mu_1, \mu_2, \mu_3) = (-4, 0, 4), (-4, 0, 6)$ and $(-6, 0, 6)$, respectively. The variance is set to $(\sigma_1, \sigma_2, \sigma_3) = (0.75, 1.5, 0.75)$. The asymptotic simulations are based on 2000 repetitions and the bootstrap simulation is based on 1000 repetitions.

Table 4: Sizes of **EM test** of $H_0 : M_0 = 2$ against $H_A : M_0 = 3$ with conditioning variables

(N, T)	(A, C, E)	(A, C, F)	(A, D, E)	(A, D, F)	(B, C, E)	(B, C, F)	(B, D, E)	(B, D, F)
(200, 2)	8.4	8.2	7.4	8.8	8.6	8.2	7.4	3.6
(500, 2)	4.6	3.2	3.2	2.2	4.8	4.8	3.6	3.6
(200, 5)	4.0	1.8	3.0	2.6	2.2	2.0	2.2	3.2
(500, 5)	2.2	1.2	1.6	1.4	3.0	2.0	1.8	2.0

Notes: A and B refer to $(\mu_1, \mu_2) = (-1, 1)$ and $(-0.5, 0.5)$, respectively, while C and D refer to $(\beta_1, \beta_2) = (1, 1)$ and $(-1, 1)$, respectively. E and F refer to $(\sigma_1, \sigma_2) = (0.3, 0.1)$ and $(0.1, 0.1)$. The mixing proportion is set to $(\alpha_1, \alpha_2) = (0.2, 0.8)$. The asymptotic simulations are based on 500 repetitions.

Table 5: Frequency of Number of Components with the Simulated Data

M	1	2	3	4	5	6	7
SHT with EM test	0	0	0	0.26	0.72	0.02	0
AIC	0	0	0	0.01	0.13	0.31	0.55
BIC	0	0	0	0.41	0.58	0.01	0

¹ The data are generated using the estimated parameters based on the Chilean textile industry with five-components and panel length $T = 3$, where $(\alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5) = (0.16076522, 0.32454077, 0.09025875, 0.35478905, 0.06964622)$, $(\mu_1, \mu_2, \mu_3, \mu_4, \mu_5) = (-1.241241, -0.33803875, 0.4480291, 0.52379553, 1.4139465)$, $(\beta_1, \beta_2, \beta_3, \beta_4, \beta_5) = (0.451833, -0.05988709, -0.2453261, -0.03106076, 0.2053708)$, $(\sigma_1, \sigma_2, \sigma_3, \sigma_4, \sigma_5) = (0.9933480, 0.4585760, 0.9954302, 0.4116855, 0.1863346)$. We use the panel length and sample size that are equal to those in the dataset, i.e., $n = 196$ and $T = 3$.

² The results are based on 100 repetitions.

³ Each cell indicates the proportion of times that the model selection indicates a M -component model.

9 Empirical Application

In this section, we conduct an empirical application of our proposed test for the number of components in a finite mixture production function model, the identification of which is analyzed in Kasahara et al. (2022). Specifically, we estimate the number of types of input elasticities in production functions using panel data from Japanese publicly traded firms in the machinery industry, as well as data from Chilean manufacturing firms.

9.1 Production Function and First Order Condition

Consider the input and output panel data of n firms over T years, $\{\{Y_{it}, V_{it}, L_{it}, K_{it}\}_{t=1}^T\}_{i=1}^N$, where Y_{it} , V_{it} , L_{it} , K_{it} represent the output, intermediate input, labor, and capital of firm i in year t , respectively. We denote the logarithm of corresponding variables by lowercase letters as $(y_{it}, v_{it}, l_{it}, k_{it})$, with, for example, $y_{it} = \log(Y_{it})$.

We employ a finite mixture specification to capture unobserved heterogeneity in a firm's input elasticities. We are interested in testing the number of production technology types. Assume there are M discrete types of production technologies and define the latent random variable $D_i \in \{1, 2, \dots, M\}$ to represent the production technology type of firm i . If $D_i = j$, then firm i is of type j . The population proportion of type j is denoted by $\alpha_j = \Pr(D = j)$. The production function for type j is Cobb-Douglas and the output is related to inputs as

$$Y_{it} = \exp(\epsilon_{it}) F_t^j(V_{it}, L_{it}, K_{it}, \omega_{it}) \quad (35)$$

with

$$F_t^j(V_{it}, L_{it}, K_{it}, \omega_{it}) := \exp(\gamma_t^j + \omega_{it}) V_{it}^{\delta_{v,j}} L_{it}^{\delta_{\ell,j}} K_{it}^{\delta_{k,j}},$$

where γ_t^j represents the aggregate productivity shock of type j in year t ; ω_{it} is the serially correlated productivity shock; and ϵ_{it} is the idiosyncratic productivity shock.

We assume that an intermediate input V_{it} is flexibly chosen by firm i after observing aggregate shock γ_t^j and serially correlated productivity shock ω_{it} . The variable ϵ_{it} represents a mean-zero i.i.d. random variable, the realization of which is unknown when the intermediate input V is selected. Denote the information available to a firm for making decisions on V_{it} by \mathcal{I}_{it} . Denote the information available to a firm for making decisions on V_{it} by \mathcal{I}_{it} .

In order to identify the intermediate input elasticity of the production function, we introduce the following assumptions (c.f., Kasahara et al. (2022)).

Assumption 5. (a) Each firm belongs to one of M types, and the probability of being type j is given by $\alpha_j = P(D_i = j)$ with $\sum_{j=1}^M \alpha_j = 1$. (b) For the j^{th} type of production technology at time t , the output is expressed in terms of input as in (35), where $\epsilon_{it} \sim N(0, \sigma_j^2)$ are i.i.d across i 's and t 's. ω_{it} follows an exogenous first-order stationary Markov process given by $\omega_{it} = h^j(\omega_{it-1}) + \eta_{it}$ where, conditional on \mathcal{I}_{it-1} , η_{it} is a mean-zero i.i.d. random variable. (c) $(\gamma_t^j, \omega_{it}) \in \mathcal{I}_{it}$ and $\epsilon_{it} \notin \mathcal{I}_{it}$.

Assumption 6. (a) Firms are price-takers in both output and input markets, where $P_{Y,t}$ and $P_{V,t}$ are the prices of output and intermediate input in year t . (b) $(P_{Y,t}, P_{V,t})$ are observed by firms at the beginning of the period before V_{it} is chosen.

Assumption 7. V_{it} 's are chosen at time t by maximizing the expected profit conditional on information \mathcal{I}_{it} at time t and conditional on the value of (K_{it}, L_{it}) . The profit maximization problem for firms with type j technology is given by

$$V_{it} = \arg \max_V P_{Y,t} \mathbb{E}[\exp(\epsilon_{it}) | D_i = j] F_t^j(V, K_{it}, L_{it}, \omega_{it}) - P_{V,t} V. \quad (36)$$

In Assumption 5(a), each firm's production function belongs to one of the M types. Assumption 5(b) assumes that the idiosyncratic productivity shock follows a normal distribution. Assumption 5(c) assumes that both the aggregate shock γ_t^j and the serially correlated productivity shock ω_{it} are observed when intermediate inputs are chosen, but idiosyncratic productivity shocks are unknown. Assumption 6 states that firms observe input and output prices when deciding on V_{it} . Assumption 7 assumes that V_{it} is chosen to maximize the current expected period profit conditional on the value of (K_{it}, L_{it}) .⁴

⁴We are agnostic about the timing of choosing K_{it} and L_{it} as long as they are either determined before V_{it} or simultaneously chosen with V_{it} . It is reasonable to assume that capital input K_{it} is determined before the value of V_{it} is

Given the above assumptions 5, 6, and 7, we derive an empirical specification based on the first-order condition of the profit maximization problem (36), following the idea developed by Gandhi et al. (2020) and extending to a finite mixture production function modeled by Kasahara et al. (2022). Note that $E[\exp(\epsilon_{it})|D_i = j] = \exp(\sigma_j^2/2)$ for $\epsilon_{it} \sim N(0, \sigma_j^2)$. Then, because $\delta_{v,j} = \frac{\partial F_i^j(V_{it}, K_{it}, L_{it})/\partial V_{it}}{F_i^j(V_{it}, K_{it}, L_{it})/V_{it}}$ for the Cobb-Douglas production function, the first-order condition with respect to V_{it} in (36) together with the production function (35) implies that:

$$s_{it} = \log \delta_{v,j} + \frac{1}{2}\sigma_j^2 - \epsilon_{it} \quad \text{for } D_i = j, \quad (37)$$

where

$$s_{it} := \log \left(\frac{P_{V,t}V_{it}}{P_{Y,t}Y_{it}} \right)$$

is the logarithm of the ratio of intermediate input cost to revenue.

Collect the observed data as $\mathbf{W}_i = \{s_{it}, \log K_{it}\}_{t=1}^T$. Let $\mu_j = \log \delta_{v,j} + \frac{1}{2}\sigma_j^2$ and define a type-specific parameter to be $\theta_j = (\mu_j, \sigma_j)$, where $\delta_{v,j}$ can be identified from θ_j as $\delta_{v,j} = \exp(\mu_j - \sigma_j^2/2)$. Collect the parameters of each type and the mixing probability as $\boldsymbol{\vartheta}_M = (\alpha_1, \dots, \alpha_{M-1}, \boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_M^\top)^\top$. Recall that $\epsilon_{it} \stackrel{iid}{\sim} N(0, \sigma_j^2)$ over i and t conditional on the technology type $D_i = j$. Then, from (37), we can write the density function of s_{i1}, \dots, s_{iT} as a mixture of type-specific likelihood density similar to the density function in equation (1):

$$f_M(\mathbf{W}_i; \boldsymbol{\vartheta}_M) = \sum_{j=1}^M \alpha_j \prod_{t=1}^T \frac{1}{\sigma_j} \phi \left(\frac{s_{it} - \mu_j}{\sigma_j} \right). \quad (38)$$

The penalized maximum likelihood estimator is defined as

$$\hat{\boldsymbol{\vartheta}}_M = \arg \max_{\boldsymbol{\vartheta}_M} \sum_{i=1}^n \log f_M(\mathbf{W}_i; \boldsymbol{\vartheta}_M) + \tilde{p}_n(\boldsymbol{\vartheta}_M).$$

As an alternative specification, we allow the elasticity of output for intermediate input to be a function of $\log K_{it}$ as $\log \delta_{v,j} = \beta_{0,j} + \beta_{k,j} \log K_{it}$. This results in the logarithm of the ratio of intermediate input cost to revenue being linearly related to $\log K_{it}$ as $s_{it} = \mu_j + \beta_{k,j} \log K_{it} - \epsilon_{it}$ for $D_i = j$ with $\mu_j = \beta_{0,j} + \frac{1}{2}\sigma_j^2$. In this case, the conditional density function of $\{s_{it}\}_{t=1}^T$ given $\{\log K_{it}\}_{t=1}^T$ is

$$f_M(\mathbf{W}_i; \boldsymbol{\vartheta}_M) = \sum_{j=1}^M \alpha_j \prod_{t=1}^T \frac{1}{\sigma_j} \phi \left(\frac{s_{it} - \mu_j - \beta_{k,j} \log K_{it}}{\sigma_j} \right). \quad (39)$$

chosen. On the other hand, labor input L_{it} may be flexibly chosen simultaneously with V_{it} after γ_t^j and ω_{it} are observed. Even when labor input is simultaneously chosen with intermediate input, equation (36) and the corresponding first-order condition characterize the intermediate input choice once we interpret L_{it} in (36) as the optimal value chosen by firm i as discussed in Akerberg et al. (2015).

In addition, we consider a specification in which we include not only $\log K_{it}$ but also $\log L_{it}$ as a regressor:

$$f_M(\mathbf{W}_i; \boldsymbol{\vartheta}_M) = \sum_{j=1}^M \alpha_j \prod_{t=1}^T \frac{1}{\sigma_j} \phi \left(\frac{s_{it} - \mu_j - \beta_{k,j} \log K_{it} - \beta_{l,j} \log L_{it}}{\sigma_j} \right). \quad (40)$$

9.2 Empirical result

We apply the EM test to two producer-level data sets to determine the number of production technology types. We used the production data from the Japanese publicly traded firms from 2003 to 2007 and the Chilean manufacturing plants from 1992 to 1996.⁵ We cleaned the data and used the firms/plants with continuous data entry for five years to ensure that we had balanced panel data. We focus on the three largest industries in terms of the number of firms and plants for each country (chemical, machine, and electronics for Japan and food products, fabricated metal products, and textiles for Chile). Table 6 presents the summary statistics for the revenue share of intermediate materials and the log of gross output in these industries. The within-industry standard deviations of the revenue share of intermediate materials are substantial across all industries, suggesting that intermediate input elasticities differ across firms within the narrowly defined industries.

⁵Please refer to Kasahara et al. (2021) and Kasahara and Rodrigue (2008) for the details of the datasets of the Japanese publicly traded firms and the Chilean manufacturing plants, respectively.

Table 6: Descriptive statistics for the revenue share of intermediate material and the log of gross output for the Japanese firms and the Chilean plants

Panel A: : Japanese publicly traded firms						
Industry	NObs	n	$\frac{P_{V,t}V_{it}}{P_{Y,t}Y_{it}}$		log(Y_{it})	
			<i>mean</i>	<i>sd</i>	<i>mean</i>	<i>sd</i>
Chemical	805	161	0.34	0.15	17.52	1.24
Machine	790	158	0.50	0.16	17.31	1.35
Electronics	775	155	0.45	0.18	17.54	1.27

Panel B: Chilean plants						
Industry	NObs	n	$\frac{P_{V,t}V_{it}}{P_{Y,t}Y_{it}}$		log(Y_{it})	
			<i>mean</i>	<i>sd</i>	<i>mean</i>	<i>sd</i>
Food products	4645	929	0.65	0.15	10.62	1.66
Fabricated metal products	1260	252	0.53	0.18	11.00	1.37
Textiles	1130	226	0.58	0.19	11.01	1.32

¹ The summary statistics are based on the Japanese firm-level data from 2003 to 2007 and the Chilean plant-level data from 1992 - 1996. All observations with $\log(V_{it}/Y_{it}) \leq -3$ and $\log(V_{it}/Y_{it}) > \log(2)$ are removed. The data set is a balanced panel, i.e., we kept firms/plants that are continuously observed for these five years.

² The variable $\frac{P_{V,t}V_{it}}{P_{Y,t}Y_{it}}$ is defined as the revenue share of the intermediate input, where $P_{V,t}$ is the average price of the intermediate input at time t , $P_{Y,t}$ is the average price of the output, V_{it} is the quantity of the intermediate input and Y_{it} is the quantity of the output.

Table 7: The EM test for Japanese producer without conditioning variables

		M=1	M=2	M=3	M=4	M=5
$T = 3$						
Chemical	<i>EM</i>	436.37***	239.83***	130.1***	126.4***	63.24***
	<i>BIC</i>	805.55	383.43	157.5	41.62	-70.46
Electronics	<i>EM</i>	563.94***	186.67***	115.82***	81.06***	47.76***
	<i>BIC</i>	814.01	264.27	91.67	-10.39	-77.2
Machine	<i>EM</i>	434.91***	194.48***	72.83***	56.94***	54.77***
	<i>BIC</i>	458.72	37.85	-142.28	-200.74	-242.71
$T = 4$						
Chemical	<i>EM</i>	629.22***	308.6***	181.39***	177.38***	96.35***
	<i>BIC</i>	1071.45	456.54	162.15	-4.99	-168.01
Electronics	<i>EM</i>	803.15***	282.32***	167.83***	106.43***	89.93***
	<i>BIC</i>	1081.48	292.68	24.54	-484.46	
Machine	<i>EM</i>	620.95***	292.52***	118.37***	102.57***	75.32***
	<i>BIC</i>	609.1	2.14	-276.04	-380.16	-467.96
$T = 5$						
Chemical	<i>EM</i>	818.38***	386.08***	219.13***	209.42***	118.25***
	<i>BIC</i>	1331.53	527.48	155.86	-48.53	-243.73
Electronics	<i>EM</i>	1024.86***	375.29***	226.01***	134.53***	126.36***
	<i>BIC</i>	1343.12	332.61	-28.32	-239.31	-359.17
Machine	<i>EM</i>	819.98***	389.69***	156.44***	149.98***	96.32***
	<i>BIC</i>	775.75	-30.17	-406.59	-548.81	-683.96

¹ The estimation is based on the revenue share of intermediate material. ² *, **, *** indicate the result is significant at 10%, 5% and 1% levels respectively.

Table 8: The EM test for Chilean producer without conditioning variables

		M=1	M=2	M=3	M=4	M=5
$T = 3$						
Food products	<i>EM</i>	805.51***	637.77***	204.92***	80.54***	72.41***
	<i>BIC</i>	422.55	-371.13	-991.96	-1176.61	-1236.82
Fabricated metal products	<i>EM</i>	238.84***	68.91***	26.24***	24.42***	21.82***
	<i>BIC</i>	719.74	496.49	444.02	433.01	425
Textiles	<i>EM</i>	229.87***	146.17***	64.76***	27.06***	29.98**
	<i>BIC</i>	635.37	418.28	288.34	236.9	223.34
$T = 4$						
Food products	<i>EM</i>	1165.08***	874.27***	257.49***	130.61***	139.59***
	<i>BIC</i>	419.47	-730.83	-1586.11	-1825.87	-1938.03
Fabricated metal products	<i>EM</i>	362.1***	120.7***	41.6***	43.68***	20.95***
	<i>BIC</i>	905.9	559.3	453.41	427.34	399.82
Textiles	<i>EM</i>	325.17***	222.28***	74.19***	47.58***	51.65***
	<i>BIC</i>	821.73	510.98	303.8	243.51	210.77
$T = 5$						
Food products	<i>EM</i>	1553.9***	1010.31***	290.02***	172.46***	155.25***
	<i>BIC</i>	471.66	-1066.71	-2057.71	-2329.38	-2484.82
Fabricated metal products	<i>EM</i>	478.94***	176.5***	58.96***	59.37***	33.19***
	<i>BIC</i>	1101.11	637.21	477.1	433.62	389.54
Textiles	<i>EM</i>	428.29***	280.46***	103.41***	56.63***	53.57***
	<i>BIC</i>	968.16	556.01	289.55	201.41	160

¹ The estimation is based on the revenue share of intermediate material. ² *, **, *** indicate the result is significant at 10%, 5% and 1% levels respectively.

To determine the number of components, we test the null hypothesis $H_0 : M = M_0$ against $H_1 : M = M_0 + 1$ by applying the EM test at the 5 percent significance level sequentially for $M_0 = 1, \dots, 5$. If we fail to reject the null hypothesis at a certain $M_0 = M$, then we conclude that there are M types of intermediate input elasticities. We consider both the models without conditioning variable (38) and the models with conditioning variable (39)-(40).

Table 7 and 8 report the result of the EM test for the model without conditioning variable (38) from the Japanese and the Chilean industries, respectively, with the panel length of $T = 3, 4, 5$ and the null model of $M = 1, \dots, 5$. For all industries in both countries and all panel lengths, we reject the null hypothesis of $H_0 : M = M_0$ for all $M_0 = 1, 2, 3, 4$, and 5 at five percent signifi-

cance level, indicating that the number of types for intermediate input elasticities is at least five types. This result reflects a considerable and persistent heterogeneity in the revenue share of intermediate materials across firms or plants, providing strong evidence for substantial heterogeneity in intermediate input elasticities across firms' production functions in Japanese and Chilean producers. Our findings serve as a caution against the conventional empirical practice of estimating the Cobb-Douglas production function, which assumes that the elasticity parameters are common across firms. Given the strong evidence of heterogeneity in production function coefficients, incorporating heterogeneity in production function coefficients in empirical applications is warranted and should be encouraged.

On the other hand, one possible reason for the estimated number of technology types being greater than five is that the assumption of the Cobb-Douglas production function may be too restrictive. When the production function is not Cobb-Douglas, the revenue share of intermediate materials generally depends on the value of production inputs (Gandhi et al., 2020). For this reason, we test the number of technology types when the revenue share of intermediate materials depends on the value of capital input as well as labor input by estimating the models (39)-(40).

Table 9 presents the results of the sequential hypothesis test and the BIC when estimating the mixture regression model with $\log K_{it}$ in (39) using data with a panel length of $T = 3$. For the Japanese Chemical, Electronics, and Machinery industries, the sequential hypothesis test suggests that the data is generated from seven to nine-component models; concurrently, the BIC selects models with at least ten components. For the Chilean Food industry, the sequential test indicates a ten-component model, while the BIC chooses an eight-component model. In contrast, the sequential hypothesis test and the BIC, respectively, select models with seven and six components for the Chilean Fabricated Metal Products industry and the Chilean Textile industry.

Table 10 reports the results for the model that includes both $\log K_{it}$ and $\log L_{it}$ as regressors. Across six industries, both the results of the sequential hypothesis test and the BIC in Table 10 select models with at least five components, providing evidence for substantial heterogeneity in production technology across firms and plants. Comparing the results of Table 10 with those of Table 9, the selected number of components for the model with $\log K_{it}$ and $\log L_{it}$ is smaller than that for the model with only $\log K$. This suggests that the number of components may be overestimated if we do not consider a sufficiently flexible production function specification by excluding some regressors.

Table 9: The EM test and the BIC (**Dependent Variable:** $\log \frac{P_{V,t}V_{it}}{P_{Y,t}Y_{it}}$, **Regressor:** $\log K_{it}$)

M_0	1	2	3	4	5	6	7	8	9	10
Japanese Chemical										
<i>EM</i>	459.4***	236.36***	125.42***	118.36***	87.63***	53.72***	38.69***	34.07**	36.47	-
<i>BIC</i>	1384.76	943.61	726.53	620.32	518.86	449.92	413.49	394.46	381.48	366.09
Japanese Electronics										
<i>EM</i>	560.06***	213.82***	116.29***	78.81***	47.05***	40.77***	27.4**	29.02	-	-
<i>BIC</i>	1332.14	788.19	593.44	495.74	434.15	406.77	385.45	372.63	367.31	351.17
Japanese Machine										
<i>EM</i>	433.19***	202.92***	80.42***	76.82***	53.83***	34.62**	55.65	-	-	-
<i>BIC</i>	1355.6	940.49	757	696.06	638.48	617.4	588.94	568.71	555.15	544.51
Chilean Food Products										
<i>EM</i>	816.06***	489.37***	169.14***	80.88***	80.63***	52.67***	31.29***	17.16**	20.55***	-60.46
<i>BIC</i>	6759.39	5962.74	5499.3	5356.47	5301.31	5241.91	5210.27	5200.71	5210.77	5222.29
Chilean Fabricated Metal Products										
<i>EM</i>	199.35***	63.25***	49.24***	30.27***	15.73**	18.25**	10.88	-	-	-
<i>BIC</i>	1923.64	1744.72	1699.97	1670.93	1661.03	1659.54	1665.08	1669.02	1680.96	1695.54
Chilean Textile										
<i>EM</i>	201.86***	95.17***	61.43***	31.17***	14.12*	17.45**	7.94	-	-	-
<i>BIC</i>	1681.91	1499.99	1424.93	1380.93	1368.65	1364.94	1365.72	1370.72	1382.83	1392.24

¹ The estimation is based on the revenue share of intermediate material using the panel data of length $T = 3$.

² *, **, *** indicate the result is significant at 10%, 5% and 1% levels respectively.

Table 10: The EM test and the BIC (**Dependent Variable:** $\log \frac{P_{V,t}V_{it}}{P_{Y,t}Y_{it}}$, **Regressors:** $\log K_{it}$ and $\log L_{it}$)

M_0	1	2	3	4	5	6	7	8	9	10
Japanese Chemical										
<i>EM</i>	412.35***	224.09***	141.59***	132.24**	121.56	-	-	-	-	-
<i>BIC</i>	1294.05	905.44	705.72	587.3	490.07	479.74	389.05	390.29	382.28	372.69
Japanese Electronics										
<i>EM</i>	573.11***	218.38***	116.07***	94.76**	47.73	-	-	-	-	-
<i>BIC</i>	1336.69	784.95	590.73	498.55	426.23	389.91	372.05	371.64	359.15	368.25
Japanese Machine										
<i>EM</i>	468.06***	204.01***	93.35***	81.62***	62.00***	37.04***	14.21	-	-	-
<i>BIC</i>	1360.56	915.69	736.2	676.26	625.66	596.45	564.34	548.7	536.78	539.64
Chilean Food Products										
<i>EM</i>	805.09***	478.64***	177.08***	84.13***	80.96***	51.97***	32.3**	19.50	-	-
<i>BIC</i>	6732.11	5952.7	5506.55	5362.27	5309.37	5257.78	5233.37	5229.9	5242.9	5258.41
Chilean Fabricated Metal Products										
<i>EM</i>	204.45***	63.57***	49.42***	28.61***	18.32	-	-	-	-	-
<i>BIC</i>	1926.06	1747.29	1709.44	1685.39	1678.71	1680.54	1685.02	1696.19	1703.21	1723.56
Chilean Textile										
<i>EM</i>	203.69***	90.69***	58.4***	32.55***	16.19	-	-	-	-	-
<i>BIC</i>	1673.99	1495.55	1431.18	1394.54	1382.59	1373.03	1368.8	1382.42	1394.3	1394.09

¹ The estimation is based on the revenue share of intermediate material using the panel data of length $T = 3$.

² *, **, *** indicate the result is significant at 10%, 5% and 1% levels respectively.

10 Conclusion

The selection of the number of components in a finite normal mixture panel regression model is a crucial practical issue that must be addressed with care. Arbitrarily choosing the number of components can result in biased estimates, invalid inference, and reduced credibility of the final outcomes. To tackle this issue, this study proposes the PLRT and an EM test and derives their asymptotic distribution for the null hypothesis of a model with M_0 components against the alternative hypothesis with $(M_0 + 1)$ components. We also develop a procedure to consistently select the number of components by sequentially applying the PLRT and EM tests. Through a simulation exercise, we demonstrate that the proposed sequential hypothesis testing procedure exhibits good performance in finite samples.

As an empirical application, we estimate the number of production technology types using producer-level panel data from Japan and Chile. We find that most industries in our dataset exhibit a level of heterogeneity that requires a five or more-component mixture model when using the Cobb-Douglas production specification or a specification in which the elasticity of inputs depends on capital and labor input linearly. This suggests strong evidence for the presence of unobserved heterogeneity in technology types. One important caveat of our empirical exercise is that the class of production functions we investigate may be restrictive. Investigating production function heterogeneity with more flexible function forms is an important future research topic.

References

- [1] Akerberg, D. A., Caves, K., and Frazer, G. (2015). Identification properties of recent production function estimators. *Econometrica*, 83(6):2411–2451.
- [2] Alexandrovich, G. (2014). A note on the article ‘inference for multivariate normal mixtures’ by j. chen and x. tan. *Journal of Multivariate Analysis*, 129:245–248.
- [3] Amengual, D., Bei, X., Carrasco, M., and Sentana, E. (2022). Score-type tests for normal mixtures. Working papers, CEMFI.
- [4] Ando, T. and Bai, J. (2016). Panel data models with grouped factor structure under unknown group membership. *Journal of Applied Econometrics*, 31(1):163–191.
- [5] Andrews, D. (1999). Estimation When a Parameter is on a Boundary. *Econometrica*, 67(6):1341–1383.
- [6] Andrews, R. L. and Currim, I. S. (2003). Retention of latent segments in regression-based marketing models. *International Journal of Research in Marketing*, 20(4):315–321.

- [7] Azaïs, J.-M., Gassiat, E., and Mercadier, C. (2009). The likelihood ratio test for general mixture models with or without structural parameter. *ESAIM: Probability and Statistics*, 13:301–327.
- [8] Balat, J., Brambilla, I., and Sasaki, Y. (2019). Heterogeneous firms: Skilled- labor productivity and the destination of exports.
- [9] Bonhomme, S. and Manresa, E. (2015). Grouped patterns of heterogeneity in panel data. *Econometrica*, 83(3):1147–1184.
- [10] Cameron, S. V. and Heckman, J. J. (1998). Life cycle schooling and dynamic selection bias: models and evidence for five cohorts of American males. *Journal of Political Economy*, 106(2):262–333.
- [11] Chen, H. and Chen, J. (2001). The likelihood ratio test for homogeneity in finite mixture models. *Canadian Journal of Statistics*, 29:201–215.
- [12] Chen, H. and Chen, J. (2003). Tests for homogeneity in normal mixtures in the presence of a structural parameter. *Statistica Sinica*, 13:351–365.
- [13] Chen, H., Chen, J., and Kalbfleisch, J. D. (2004). Testing for a finite mixture model with two components. *Journal of the Royal Statistical Society, Series B*, 66:95–115.
- [14] Chen, J. (1995). Optimal rate of convergence for finite mixture models. *Annals of Statistics*, 23(1):221–233.
- [15] Chen, J. and Li, P. (2009). Hypothesis test for normal mixture models: The EM approach. *Annals of Statistics*, 37:2523–2542.
- [16] Chen, J. and Tan, X. (2009). Inference for multivariate normal mixtures. *Journal of Multivariate Analysis*, 100(7):1367–1383.
- [17] Chen, X., Ponomareva, M., and Tamer, E. (2014). Likelihood inference in some finite mixture models. *Journal of Econometrics*, 182(1):87–99.
- [18] Chernoff, H. and Lander, E. (1995). Asymptotic distribution of the likelihood ratio test that a mixture of two binomials is a single binomial. *Journal of Statistical Planning and Inference*, 43:19–40.
- [19] Dacunha-Castelle, D. and Gassiat, E. (1999). Testing the order of a model using locally conic parametrization: Population mixtures and stationary ARMA processes. *Annals of Statistics*, 27:1178–1209.

- [20] Deb, P. and Trivedi, P. K. (1997). Demand for medical care by the elderly: a finite mixture approach. *Journal of Applied Econometrics*, 12(3):313–336.
- [21] Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22.
- [22] Doraszelski, U. and Jaumandreu, J. (2018). Measuring the bias of technological change. *Journal of Political Economy*, 126(3):1027–1084.
- [23] Foutz, R. V. and Srivastava, R. C. (1977). The performance of the likelihood ratio test when the model is incorrect. *The Annals of Statistics*, 5(6):1183–1194.
- [24] Gandhi, A., Navarro, S., and Rivers, D. A. (2020). On the Identification of Gross Output Production Functions. *Journal of Political Economy*, 128(8):2973–3016.
- [25] Garel, B. (2001). Likelihood ratio test for univariate Gaussian mixture. *Journal of Statistical Planning and Inference*, 96:325–350.
- [26] Garel, B. (2005). Asymptotic theory of the likelihood ratio test for the identification of a mixture. *Journal of Statistical Planning and Inference*, 131:271–296.
- [27] Ghosh, J. K. and Sen, P. K. (1985). On the asymptotic performance of the log-likelihood ratio statistic for the mixture model and related results. In Le Cam, L. and Olshen, R., editors, *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, volume 2, pages 789–806. Belmont, CA: Wadsworth.
- [28] Hao, J. (2017). *NormalRegPanelMixture: Finite Mixture Model with Normal Panel Data*. R package version 1.0.
- [29] Hartigan, J. (1985). Failure of log-likelihood ratio test. In *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, volume 2, pages 807–810. University of California Press 2. Berkeley.
- [30] Heckman, J. and Singer, B. (1984). A method for minimizing the impact of distributional assumptions in econometric models for duration data. *Econometrica*, 52(2):271–320.
- [31] Kamakura, W. and Russell, G. (1989). A probabilistic choice model for market segmentation and elasticity structure. *Journal of Marketing Research*, 26(4):379–390.
- [32] Kasahara, H. and Rodrigue, J. (2008). Does the use of imported intermediates increase productivity? plant-level evidence. *Journal of Development Economics*, 87(1):106–118.

- [33] Kasahara, H., Schrimpf, P., and Suzuki, M. (2022). Identification and estimation of production function with unobserved heterogeneity. Technical report, ESRI Discussion Paper Series No.368.
- [34] Kasahara, H. and Shimotsu, K. (2009). Nonparametric Identification of Finite Mixture Models of Dynamic Discrete Choices. *Econometrica*, 77(1):135–175.
- [35] Kasahara, H. and Shimotsu, K. (2012). Testing the number of components in finite mixture models.
- [36] Kasahara, H. and Shimotsu, K. (2014). Non-parametric identification and estimation of the number of components in multivariate mixtures. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 76(1):97–111.
- [37] Kasahara, H. and Shimotsu, K. (2015). Testing the number of components in normal mixture regression models. *Journal of the American Statistical Association*, 110(512):1632–1645.
- [38] Kasahara, H. and Shimotsu, K. (2019). Testing the Order of Multivariate Normal Mixture Models.
- [39] Kasahara, H., Suzuki, M., and Sawada, Y. (2021). The effect of bank recapitalization policy on credit allocation and corporate investment: Evidence from a banking crisis in japan: Economic and social research institute. Technical report, ESRI Discussion Paper Series No.365.
- [40] Keane, M. P. and Wolpin, K. I. (1997). The career decisions of young men. *Journal of Political Economy*, 105(3):473–522.
- [41] Lehmann, E. L. and Romano, J. P. (2005). *Testing Statistical Hypotheses*. Springer, third edition edition.
- [42] Lemdani, M. and Pons, O. (1997). Likelihood ratio tests for genetic linkage. *Statistics and Probability Letters*, 33:15–22.
- [43] Levinsohn, J. and Petrin, A. (2003). Estimating Production Functions Using Inputs to Control for Unobservables. *Review of Economic Studies*, pages 317–341.
- [44] Li, P. and Chen, J. (2010). Testing the order of a finite mixture. *Journal of the American Statistical Association*, 105:1084–1092.
- [45] Li, P., Chen, J., and Marriott, P. (2009). Non-finite Fisher information and homogeneity: An EM approach. *Biometrika*, 96:411–426.
- [46] Li, T. and Sasaki, Y. (2017). Constructive identification of heterogeneous elasticities in the cobb-douglas production function.

- [47] Lin, C.-C. and Ng, S. (2012). Estimation of panel data models with parameter heterogeneity when group membership is unknown. *Journal of Econometric Methods*, 1:42–55.
- [48] Lindsay, B. G. (1995). Mixture Models: Theory, Geometry and Applications NSF-CBMS Regional Conference Series in Probability and Statistics. Source: NSF-CBMS Regional Conference Series in Probability and Statistics, 5:1–163.
- [49] Liu, G., Fu, Y., Li, P., and Pu, X. (2018). Using differential variability to increase the power of the homogeneity test in a two-sample problem. *Statistica Sinica*, 28(1):27–41.
- [50] Liu, X. and Shao, Y. (2003). Asymptotics for likelihood ratio tests under loss of identifiability. *Annals of Statistics*, 31:807–832.
- [51] Lu, X. and Su, L. (2017). Determining the number of groups in latent panel structures with an application to income and democracy. *Quantitative Economics*, 8(3):729–760.
- [52] McLachlan, G. and Peel, D. (2004). *Finite Mixture Models*.
- [53] Melnykov, V. and Maitra, R. (2010). Finite mixture models and model-based clustering. *Statistics Survey*, 4:80–116.
- [54] Newey, W. K. and McFadden, D. (1994). Large sample estimation and hypothesis testing. In *Handbook of Econometrics*, volume 4, pages 2111–2245. Elsevier.
- [55] Niu, X., Li, P., and Zhang, P. (2011). Testing homogeneity in a multivariate mixture model. *The Canadian Journal of Statistics / La Revue Canadienne de Statistique*, 39(2):218–238.
- [56] Olley, G. S. and Pakes, A. (1996). The Dynamics of Productivity in the Telecommunications Equipment Industry. *Econometrica*, pages 1263–1297.
- [57] Robin, J.-M. and Smith, R. J. (2000). Tests of rank. *Econometric Theory*, 16(2):151–175.
- [58] Shen, J. and He, X. (2015). Inference for subgroup analysis with a structured logistic-normal mixture model. *Journal of the American Statistical Association*, 110(509):303–312.
- [59] Su, L., Zhentao, S., and Phillips, P. (2016). Identifying latent structures in panel data. *Econometrica*, 84:2215–2264.
- [60] Titterton, D. M., Smith, A. F., and Makov, U. E. (1985). *Statistical Analysis of Finite Mixture Distributions*. Wiley.
- [61] White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica*, 50(1):1–25.
- [62] Zhu, H.-T. and Zhang, H. (2004). Hypothesis testing in mixture regression models. *Journal of the Royal Statistical Society, Series B*, 66:3–16.

A Proofs of propositions

Proof of Proposition 1. We first consider a model with intercept parameter and variance parameter but without covariates with $\mathbf{W}_i = \{y_{it}\}_{t=1}^T$.

Define

$$s_i^2 = \frac{1}{T-1} \sum_{t=1}^T (Y_{it} - \bar{Y}_i)^2 \quad \text{with} \quad \bar{Y}_i = \frac{1}{T} \sum_{t=1}^T Y_{it},$$

where s_i^2 follows the chi-square distribution with the $T-1$ degrees of freedom. Let $i^* = \arg \min_{i=1, \dots, n} \{s_i^2\}$ so that $s_{i^*}^2 = \min\{s_1^2, \dots, s_n^2\}$ be the minimum of s_i^2 across all i 's. We consider a sequence of parameters $\boldsymbol{\vartheta}_{2,n} = (\alpha_n, \boldsymbol{\theta}_{1,n}^\top, \boldsymbol{\theta}_{2,n}^\top)^\top$ with $\alpha_n = 1/n$, $\boldsymbol{\theta}_{1,n} = (\mu_{1,n}, \sigma_{1,n}^2)^\top = (\bar{Y}_{i^*}, s_{i^*}^2)^\top$, and $\boldsymbol{\theta}_{2,n} = \boldsymbol{\theta}^* = (\mu^*, \sigma^*)^\top$ for all n . Because $LR_n^*(\boldsymbol{\vartheta}_{2,n}) \leq LR_n^*(\tilde{\boldsymbol{\vartheta}}_{2,n})$, it suffices to show that $LR_n^*(\boldsymbol{\vartheta}_{2,n})$ is unbounded in probability.

Define

$$\ell(\mathbf{W}_i; \boldsymbol{\theta}) := \log f(\mathbf{W}_i; \boldsymbol{\theta}) = -\frac{T}{2} \log \sigma^2 - \frac{T}{2} \log(2\pi) - \frac{1}{2} \sum_{t=1}^T \left(\frac{Y_{it} - \mu}{\sigma} \right)^2.$$

Then, the likelihood ratio test statistic for a two-component mixture is written as:

$$\begin{aligned} LR_n^*(\boldsymbol{\vartheta}_{2,n}) &= 2 \left\{ \sum_{i=1}^n \log \left(\alpha_n \prod_{t=1}^T \frac{1}{\sigma_{1,n}} \phi \left(\frac{Y_{it} - \mu_{1,n}}{\sigma_{1,n}} \right) + (1 - \alpha_n) \prod_{t=1}^T \frac{1}{\sigma^*} \phi \left(\frac{Y_{it} - \mu^*}{\sigma^*} \right) \right) - \sum_{i=1}^n \ell(\mathbf{W}_i; \boldsymbol{\theta}^*) \right\} \\ &= 2 \sum_{i \neq i^*} \left\{ \log \left(\exp(\log \alpha_n + \ell(\mathbf{W}_i; \boldsymbol{\theta}_{1,n})) + \exp(\log(1 - \alpha_n) + \ell(\mathbf{W}_i; \boldsymbol{\theta}^*)) \right) - \ell(\mathbf{W}_i; \boldsymbol{\theta}^*) \right\} \\ &\quad + 2 \left\{ \log \left(\exp(\log \alpha_n + \ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}_{1,n})) + \exp(\log(1 - \alpha_n) + \ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}^*)) \right) - \ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}^*) \right\}. \end{aligned} \quad (41)$$

The first term on the right hand side of (41) can be re-written as:

$$= 2(n-1) \log \left(\frac{n-1}{n} \right) + 2 \sum_{i \neq i^*} \log \left(1 + \frac{1}{n-1} \exp(\ell(\mathbf{W}_i; \boldsymbol{\theta}_{1,n}) - \ell(\mathbf{W}_i; \boldsymbol{\theta}^*)) \right),$$

which is bounded from below by -1 as $n \rightarrow \infty$ because $\lim_{n \rightarrow \infty} 2(n-1) \log \left(\frac{n-1}{n} \right) = -1$ and $\log \left(1 + \frac{1}{n-1} \exp(\ell(\mathbf{W}_i; \boldsymbol{\theta}_{1,n}) - \ell(\mathbf{W}_i; \boldsymbol{\theta}^*)) \right) \geq 0$ for all n .

The second term on the right-hand side of (41) is written as

$$2\{-\log n + \ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}_{1,n})\} + 2 \log \left(1 + (n-1) \exp(\ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}^*) - \ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}_{1,n})) \right) - 2\ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}^*), \quad (42)$$

where $2\{-\log n + \ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}_{1,n})\}$ diverges to infinity as $n \rightarrow \infty$ by Lemma 1, while the second

term in (42) is bounded below from zero and the third term, is bounded in probability because $\ell(\mathbf{W}_{i^*}; \boldsymbol{\theta}^*) = O_p(1)$. Therefore, for any $M < \infty$, we have $\Pr\left(LR_n^*(\boldsymbol{\vartheta}_{2,n}) \leq M\right) \rightarrow 0$ as $n \rightarrow \infty$. The stated result follows from $LR_n^*(\boldsymbol{\vartheta}_{2,n}) \leq LR_n^*(\hat{\boldsymbol{\vartheta}}_{2,n})$ for all n .

For a model with covariates, we may consider a sequence of parameters $\boldsymbol{\vartheta}_{2,n} = (\alpha_n, \boldsymbol{\theta}_{1,n}^\top, \boldsymbol{\theta}_{2,n}^\top, \boldsymbol{\gamma}_n^\top)^\top$ with $\alpha_n = 1/n$, $\boldsymbol{\theta}_{1,n} = (\mu_{1,n}, \sigma_{1,n}^2, \boldsymbol{\beta}_{1,n}^\top)^\top = (\bar{Y}_{i^*} - \bar{\mathbf{Z}}_{i^*}^\top \boldsymbol{\gamma}^*, s_{i^*}^2, \mathbf{0}^\top)^\top$ with $\bar{\mathbf{Z}}_{i^*} = (1/T) \sum_{t=1}^T \mathbf{Z}_{it}$, $\boldsymbol{\theta}_{2,n} = \boldsymbol{\theta}^* = (\mu^*, \sigma^*, (\boldsymbol{\beta}^*)^\top)^\top$, and $\boldsymbol{\gamma}_n = \boldsymbol{\gamma}^*$. Then, repeating the above argument, the state result follows. \square

Proof of Proposition 2. The stated result follows from repeating the proof of Proposition 6. \square

Proof of Proposition 3. The proof follows that of Proposition 2 in Kasahara and Shimotsu (2012). For a vector \mathbf{x} and a function $f(\mathbf{x})$, let $\nabla_{\mathbf{x}^k} f(\mathbf{x})$ denote its k -th derivative with respect to \mathbf{x} , which can be a multidimensional array. Observe that, for any finite k and for a neighborhood \mathcal{N} of $\boldsymbol{\psi}^*$, we obtain

$$\begin{aligned} E\|\nabla_{\boldsymbol{\psi}^k} g(\mathbf{W}_i; \boldsymbol{\psi}^*, \alpha)/g(\mathbf{W}_i; \boldsymbol{\psi}^*, \alpha)\|^2 &< \infty, \\ E\|\sup_{\boldsymbol{\psi} \in \Theta_{\boldsymbol{\psi}} \cap \mathcal{N}} \nabla_{\boldsymbol{\psi}^k} \log g(\mathbf{W}_i; \boldsymbol{\psi}, \alpha)\|^2 &< \infty, \end{aligned} \quad (43)$$

because each element of $\nabla_{\boldsymbol{\psi}^k} \log g(y|\mathbf{x}, \mathbf{z}; \boldsymbol{\psi}, \alpha)$ is written as a sum of products of Hermite polynomials. Note also that the following holds:

$$\nabla_{\eta\lambda_j} L_n(\boldsymbol{\psi}^*, \alpha) = 0, \quad \nabla_{\lambda_i\lambda_j\lambda_k} L_n(\boldsymbol{\psi}^*, \alpha) = O_p(n^{1/2}), \quad (44)$$

$$\nabla_{\eta\eta\lambda_i} L_n(\boldsymbol{\psi}^*, \alpha) = O_p(n), \quad \nabla_{\eta\eta\eta} L_n(\boldsymbol{\psi}^*, \alpha) = O_p(n), \quad (45)$$

where equation (44) follows from Proposition 3(a)(c) and (43) while equation (45) is a simple consequence of (43). Furthermore, for a neighborhood \mathcal{N} of $\boldsymbol{\psi}^*$,

$$\sup_{\boldsymbol{\psi} \in \Theta_{\boldsymbol{\psi}} \cap \mathcal{N}} \left| n^{-1} \nabla^{(4)} L_n(\boldsymbol{\psi}, \alpha) - E \nabla^{(4)} \log g(\mathbf{W}_i; \boldsymbol{\psi}, \alpha) \right| = o_p(1), \quad (46)$$

$$E \nabla^{(4)} g(\mathbf{W}_i; \boldsymbol{\psi}, \alpha) \text{ is continuous in } \boldsymbol{\psi} \in \Theta_{\boldsymbol{\psi}} \cap \mathcal{N}. \quad (47)$$

Equations (46) and (47) follow from Lemma 2.4 of Newey and McFadden (1994) and the fact that $\nabla_{\boldsymbol{\psi}^k} \log g(\mathbf{w}; \boldsymbol{\psi}, \alpha)$ is written as a sum of products of Hermite polynomials.

Taking a fourth-order Taylor expansion of $L_n(\boldsymbol{\psi}, \alpha)$ around $\boldsymbol{\psi}^*$ and using (43) and (44), we can

write $L_n(\boldsymbol{\psi}, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha)$ as the sum of relevant terms and the remainder term as follows:

$$L_n(\boldsymbol{\psi}, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha) = \nabla_{\boldsymbol{\eta}} L_n^*(\boldsymbol{\eta} - \boldsymbol{\eta}^*) + \frac{1}{2!} (\boldsymbol{\eta} - \boldsymbol{\eta}^*)^\top \nabla_{\boldsymbol{\eta} \boldsymbol{\eta}^\top} L_n^*(\boldsymbol{\eta} - \boldsymbol{\eta}^*) + \frac{1}{2!} \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \nabla_{\lambda_i \lambda_j} L_n^* \lambda_i \lambda_j \quad (48)$$

$$+ \frac{3}{3!} \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} (\boldsymbol{\eta} - \boldsymbol{\eta}^*)^\top \nabla_{\boldsymbol{\eta} \lambda_i \lambda_j} L_n^* \lambda_i \lambda_j \quad (49)$$

$$+ \frac{1}{4!} \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \sum_{\ell=1}^{q+2} \nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} L_n^* \lambda_i \lambda_j \lambda_k \lambda_\ell + R_n(\boldsymbol{\psi}, \alpha), \quad (50)$$

where ∇L_n^* denotes the derivative of $L_n(\boldsymbol{\psi}, \alpha)$ evaluated at $(\boldsymbol{\psi}^*, \alpha)$. In view of (44)-(45), the remainder term is written as

$$R_n(\boldsymbol{\psi}, \alpha) = O_p(n^{1/2}) \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \lambda_i \lambda_j \lambda_k + O_p(n) \left(\sum_{i=1}^{q+2} \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^2 \lambda_i + \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^3 \right) \quad (51)$$

$$+ O_p(n) \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} (\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^4 + \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^3 |\lambda_i| + \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^2 |\lambda_i \lambda_j| + \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| |\lambda_i \lambda_j \lambda_k|) \quad (52)$$

$$+ \frac{1}{4!} \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \sum_{\ell=1}^{q+2} \{\nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} L_n(\boldsymbol{\psi}^\dagger, \alpha) - \nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} L_n(\boldsymbol{\psi}^*, \alpha)\} \lambda_i \lambda_j \lambda_k \lambda_\ell \quad (53)$$

with $\boldsymbol{\psi}^\dagger$ being between $\boldsymbol{\psi}$ and $\boldsymbol{\psi}^*$. Because $\|\sqrt{nt}(\boldsymbol{\psi}, \alpha)\|^2 = n\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^2 + n \sum_{i=1}^{q+2} \sum_{j=1}^i \alpha^2 (1 - \alpha)^2 |\lambda_i \lambda_j|^2$, the right hand side of (51) and the terms in (52) are bounded by $O_p(1)(\|\sqrt{nt}(\boldsymbol{\psi}, \alpha)\| + \|\sqrt{nt}(\boldsymbol{\psi}, \alpha)\|^2)(\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| + \|\lambda\|)$. In view of (46) and (47), (53) is bounded by $\|\sqrt{nt}(\boldsymbol{\psi}, \alpha)\|^2 [d(\boldsymbol{\psi}^\dagger) + o_p(1)]$ with $d(\boldsymbol{\psi}^\dagger) \rightarrow 0$ as $\boldsymbol{\psi}^\dagger \rightarrow \boldsymbol{\psi}^*$, where a function $d(\boldsymbol{\psi}^\dagger)$ corresponds to $n^{-1} \mathbb{E}[\nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} L_n(\boldsymbol{\psi}^\dagger, \alpha) - \nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} L_n(\boldsymbol{\psi}^*, \alpha)]$. Therefore, $R_n(\boldsymbol{\psi}, \alpha) = (1 + \|\sqrt{nt}(\boldsymbol{\psi}, \alpha)\|)^2 [d(\boldsymbol{\psi}^\dagger) + o_p(1) + O_p(\|\boldsymbol{\psi} - \boldsymbol{\psi}^*\|)]$, and part (a) follows.

Part (b) follows from Lemman 3(c)(d), the Lindeberg-Levy central limit theorem, and the finiteness of \mathcal{I} in part (c).

For part (c), we first provide the formula of \mathcal{I}_n . Partition \mathcal{I}_n as

$$\mathcal{I}_n = \begin{pmatrix} \mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\eta}^n} & \mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\lambda}^n} \\ \mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\lambda}^n}^\top & \mathcal{I}_{\boldsymbol{\lambda}^n} \end{pmatrix}, \quad \mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\eta}^n} : (p+q+2) \times (p+q+2), \quad \mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\lambda}^n} : (p+q+2) \times q_\lambda, \quad \mathcal{I}_{\boldsymbol{\lambda}^n} : q_\lambda \times q_\lambda,$$

where q_λ represents the number of unique terms in $\sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \sum_{\ell=1}^{q+2} \lambda_i \lambda_j \lambda_k \lambda_\ell$. $\mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\eta}^n}$ is given by $\mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\eta}^n} = -n^{-1} \nabla_{\boldsymbol{\eta} \boldsymbol{\eta}^\top} L_n(\boldsymbol{\psi}^*, \alpha)$. For $\mathcal{I}_{\boldsymbol{\eta} \boldsymbol{\lambda}^n}$, let $A_{ij} = n^{-1} \nabla_{\boldsymbol{\eta} \lambda_i \lambda_j} L_n(\boldsymbol{\psi}^*, \alpha)$

and write the term in (49) as $(n/2) \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} (\boldsymbol{\eta} - \boldsymbol{\eta}^*)^\top A_{ij} \lambda_i \lambda_j = n \sum_{i=1}^{q+2} \sum_{j=1}^i c_{ij} (\boldsymbol{\eta} - \boldsymbol{\eta}^*)^\top A_{ij} \lambda_i \lambda_j$, where the c_{ij} 's are defined when we introduce $\tilde{\nabla}_{\boldsymbol{\theta}\boldsymbol{\theta}^\top} f^*$ after (9). Then, by defining $\mathcal{I}_{\boldsymbol{\eta}\boldsymbol{\lambda}n} = -(c_{11}A_{11}, \dots, c_{qq}A_{qq}, c_{12}A_{12}, \dots, c_{q-1,q}A_{q-1,q})/\alpha(1-\alpha)$, the term in (49) equals $-n(\boldsymbol{\eta} - \boldsymbol{\eta}^*)^\top \mathcal{I}_{\boldsymbol{\eta}\boldsymbol{\lambda}n} [\alpha(1-\alpha)v(\boldsymbol{\lambda})]$. For $\mathcal{I}_{\boldsymbol{\lambda}n}$, define $B_{ijkl} = n^{-1}(8/4!) \nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} L_n(\boldsymbol{\psi}^*, \alpha)$ so that the first term in (50) is written as $(n/8) \sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \sum_{\ell=1}^{q+2} B_{ijkl} \lambda_i \lambda_j \lambda_k \lambda_\ell = (n/2) \sum_{i=1}^{q+2} \sum_{j=1}^i \sum_{k=1}^{q+2} \sum_{\ell=1}^k c_{ij} c_{k\ell} B_{ijkl} \lambda_i \lambda_j \lambda_k \lambda_\ell$. Define $\mathcal{I}_{\boldsymbol{\lambda}n}$ such that the $(ij, k\ell)$'s element of $\mathcal{I}_{\boldsymbol{\lambda}n}$ is $-c_{ij} c_{k\ell} B_{ijkl} / \alpha^2 (1-\alpha)^2$, where the ij 's run over $\{(1, 1), \dots, (q, q), (1, 2), \dots, (q-1, q)\}$. Then, the first term in (50) equals $-(n/2) [\alpha(1-\alpha)v(\boldsymbol{\lambda})]' \mathcal{I}_{\boldsymbol{\lambda}n} [\alpha(1-\alpha)v(\boldsymbol{\lambda})]$. With this definition of \mathcal{I}_n , the expansion (48)-(50) is written as (12) in terms of $\sqrt{nt}(\boldsymbol{\psi}, \alpha)$.

We now show that $\mathcal{I}_n \rightarrow_p \mathcal{I}$. $\mathcal{I}_{\boldsymbol{\eta}n} \rightarrow_p \mathcal{I}_{\boldsymbol{\eta}}$ holds trivially. For $\mathcal{I}_{\boldsymbol{\eta}\boldsymbol{\lambda}n}$, it follows from Lemma 3(c) and the law of large numbers that $A_{ij} \rightarrow_p -\mathbb{E}[\nabla_{\boldsymbol{\eta}} l(\mathbf{W}; \boldsymbol{\psi}^*, \alpha) \nabla_{\lambda_i \lambda_j} l(\mathbf{W}; \boldsymbol{\psi}^*, \alpha)]$, giving $\mathcal{I}_{\boldsymbol{\eta}\boldsymbol{\lambda}n} \rightarrow_p E[\mathbf{s}_{\boldsymbol{\eta}} \mathbf{s}_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^\top / \alpha(1-\alpha)] = \mathcal{I}_{\boldsymbol{\eta}\boldsymbol{\lambda}}$. For $\mathcal{I}_{\boldsymbol{\lambda}n}$, Lemma 3(d) and the law of large numbers imply that $\sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \sum_{\ell=1}^{q+2} B_{ijkl} \lambda_i \lambda_j \lambda_k \lambda_\ell \rightarrow_p -\sum_{i=1}^{q+2} \sum_{j=1}^{q+2} \sum_{k=1}^{q+2} \sum_{\ell=1}^{q+2} E[\nabla_{\lambda_i \lambda_j} l(\mathbf{W}; \boldsymbol{\psi}^*, \alpha) \nabla_{\lambda_k \lambda_\ell} l(\mathbf{W}; \boldsymbol{\psi}^*, \alpha)] \lambda_i \lambda_j \lambda_k \lambda_\ell$, where the factor $(8/4!) = 1/3$ in B_{ijkl} and the three derivatives on the right hand side of Lemma 3(d) cancel each other. Therefore, we have $\mathcal{I}_{\boldsymbol{\lambda}n} \rightarrow_p E[\mathbf{s}_{\boldsymbol{\lambda}\boldsymbol{\lambda}} \mathbf{s}_{\boldsymbol{\lambda}\boldsymbol{\lambda}}^\top / \alpha^2 (1-\alpha)^2] = \mathcal{I}_{\boldsymbol{\lambda}}$, and $\mathcal{I}_n \rightarrow_p \mathcal{I}$ follows.

We complete the proof of part (c) by showing that $\mathcal{I} = E[\mathbf{s}(\mathbf{W})\mathbf{s}(\mathbf{W})^\top]$ is finite and non-singular. Note that $\mathbf{s}(\mathbf{W})$ can be expressed in Hermite polynomials as in (74). Then, the finiteness of \mathcal{I} follows from Assumption 2(a) and the definition of Hermite polynomials.

To show that \mathcal{I} is positive definite, it suffices to show that there exists no multi-collinearity in $\mathbf{s}(\mathbf{w})$. Suppose, on the contrary, that $\mathbf{s}(\mathbf{w})$ is multi-collinear and that there exists a non-zero vector \mathbf{a} that solves the equation $\mathbf{a}^\top \mathbf{s}(\mathbf{w}) = 0$ for all values of \mathbf{w} . Partition $\mathbf{s}(\mathbf{w})$ as $\mathbf{s}(\mathbf{w}) = (\mathbf{s}_{(\mu)}^\top, \mathbf{s}_{(\beta)}^\top)^\top$ with $\mathbf{s}_{(\mu)} = (s_\mu, s_\sigma, s_{\lambda\mu\mu}, s_{\lambda\mu\sigma}, s_{\lambda\sigma\sigma})^\top$ and $\mathbf{s}_{(\beta)} = (\mathbf{s}_{\beta}^\top, \mathbf{s}_{\gamma}^\top, \mathbf{s}_{\mu\beta}^\top, \mathbf{s}_{\sigma\beta}^\top, \mathbf{s}_{\lambda\beta\beta}^\top)^\top$, where $\mathbf{s}(\mathbf{w})$ is defined in (9) and (74). Similarly, partition \mathbf{a} as $\mathbf{a} = (\mathbf{a}_{(\mu)}^\top, \mathbf{a}_{(\beta)}^\top)^\top$ so that

$$\mathbf{a}^\top \mathbf{s}(\mathbf{w}) = \mathbf{a}_{(\mu)}^\top \mathbf{s}_{(\mu)} + \mathbf{a}_{(\beta)}^\top \mathbf{s}_{(\beta)}. \quad (54)$$

By Assumption 2(b) and the property of Hermite polynomials, if $\mathbf{a}^\top \mathbf{s}(\mathbf{w}) = \mathbf{0}$ for all \mathbf{w} , then $\mathbf{a}_{(\beta)} = \mathbf{0}$.

Then, in view of (54), the stated result follows if we can show that $\mathbf{a}_{(\mu)}^\top \mathbf{s}_{(\mu)} = \mathbf{0}$ for all \mathbf{w} implies $\mathbf{a}_{(\mu)} = \mathbf{0}$. Suppose that

$$\begin{aligned} \mathbf{a}_{(\mu)}^\top \mathbf{s}_{(\mu)} &= a_\mu \sum_{t=1}^T H_{i,t}^{1*} + (a_\sigma + a_{\lambda\mu\mu}) \sum_{t=1}^T H_{i,t}^{2*} + \frac{a_{\lambda\mu\mu}}{2} \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} H_{i,s}^{1*} \\ &+ a_{\lambda\mu\sigma} \sum_{t=1}^T H_{i,t}^{3*} + a_{\lambda\mu\sigma} \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} H_{i,s}^{2*} + 3a_{\lambda\sigma\sigma} \sum_{t=1}^T H_{i,t}^{4*} + \frac{a_{\lambda\sigma\sigma}}{2} \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{2*} H_{i,s}^{2*} = 0 \end{aligned}$$

for all \mathbf{w} , where $H_{i,t}^{j*}$ for $j = 1, 2, 3$ is defined in (73) in Appendix B.2.

Because the above equation hold for all values of \mathbf{w} , with the property of the Hermite polynomials, we have $a_\mu = 0, (a_\sigma + a_{\lambda\mu\mu}) = 0, a_{\lambda\mu\mu} = 0, a_{\lambda\mu\sigma} = 0, a_{\lambda\sigma\sigma} = 0$. This implies that $\mathbf{a}_{(\mu)} = 0$. Therefore, there exists no multi-collinearity in $s(\mathbf{w})$, and \mathcal{I} is non-singular, proving part (c). \square

Proof of Proposition 4. The proof is similar to that of Proposition 3 in Kasahara and Shimotsu (2015).

The proof of part (a) closely follows the proof of Theorem 1 of Andrews (1999). Let $\mathbf{T}_n := \mathcal{I}_n^{1/2} \sqrt{nt}(\hat{\boldsymbol{\psi}}_\alpha, \alpha)$. Then, in view of (12), we have

$$\begin{aligned} o_p(1) &\leq L_n(\hat{\boldsymbol{\psi}}_\alpha, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha) \\ &= \mathbf{T}'_n \mathcal{I}_n^{-1/2} \mathbf{S}_n - \frac{1}{2} \|\mathbf{T}_n\|^2 + R_n(\hat{\boldsymbol{\psi}}_\alpha, \alpha) \\ &= O_p(\|\mathbf{T}_n\|) - \frac{1}{2} \|\mathbf{T}_n\|^2 + (1 + \|\mathcal{I}_n^{-1/2} \mathbf{T}_n\|)^2 o_p(1) \\ &= \|\mathbf{T}_n\| O_p(1) - \frac{1}{2} \|\mathbf{T}_n\|^2 + o_p(\|\mathbf{T}_n\|) + o_p(\|\mathbf{T}_n\|^2) + o_p(1), \end{aligned}$$

where the third equality holds because $\mathcal{I}_n^{-1/2} \mathbf{S}_n = O_p(1)$ and $R_n(\hat{\boldsymbol{\psi}}_\alpha, \alpha) = o_p((1 + \|\mathcal{I}_n^{-1/2} \mathbf{T}_n\|)^2)$ from Propositions 2 and 3. Rearranging this equation yields $\|\mathbf{T}_n\|^2 \leq 2\|\mathbf{T}_n\| O_p(1) + o_p(1)$. Denote the $O_p(1)$ term by ς_n . Then, $(\|\mathbf{T}_n\| - \varsigma_n)^2 \leq \varsigma_n^2 + o_p(1) = O_p(1)$; taking its square root gives $\|\mathbf{T}_n\| \leq O_p(1)$. In conjunction with $\mathcal{I}_n \rightarrow_p \mathcal{I}$, we obtain $\sqrt{nt}(\hat{\boldsymbol{\psi}}_\alpha, \alpha) = O_p(1)$, and part (a) follows.

For part (b), noting that $L_n(\boldsymbol{\psi}^*, \alpha) = L_{0,n}(\boldsymbol{\gamma}_0^*, \boldsymbol{\theta}_0^*)$, write

$$LR_n = \max_{\alpha \in [\epsilon, 1-\epsilon]} 2\{L_n(\hat{\boldsymbol{\psi}}_\alpha, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha)\} - 2\{L_{0,n}(\hat{\boldsymbol{\gamma}}_0, \hat{\boldsymbol{\theta}}_0) - L_{0,n}(\boldsymbol{\gamma}_0^*, \boldsymbol{\theta}_0^*)\}. \quad (55)$$

Define

$$\mathbf{S}_n = \begin{pmatrix} \mathbf{S}_{\eta n} \\ \mathbf{S}_{\lambda n} \end{pmatrix} := \begin{pmatrix} n^{-1/2} \sum_{i=1}^n \mathbf{s}_\eta(\mathbf{W}_i) \\ n^{-1/2} \sum_{i=1}^n \mathbf{s}_{\lambda\lambda}(\mathbf{W}_i) \end{pmatrix}, \quad \begin{aligned} \mathbf{S}_{\lambda, \eta n} &:= \mathbf{S}_{\lambda n} - \mathcal{I}_{\lambda\eta} \mathcal{I}_\eta^{-1} \mathbf{S}_{\eta n}, & \mathbf{G}_{\lambda, \eta n} &:= \mathcal{I}_{\lambda, \eta}^{-1} \mathbf{S}_{\lambda, \eta n}, \\ \mathbf{t}_{\eta, \lambda} &:= \mathbf{t}_\eta - \mathcal{I}_\eta \mathcal{I}_{\eta\lambda}^{-1} \mathbf{t}_\lambda(\boldsymbol{\lambda}, \alpha), \end{aligned}$$

and split the quadratic form in (12) to obtain

$$2\{L_n(\boldsymbol{\psi}, \alpha) - L_n(\boldsymbol{\psi}^*, \alpha)\} = B_n(\sqrt{nt} \mathbf{t}_{\eta, \lambda}) + C_n(\sqrt{nt} \mathbf{t}_\lambda(\boldsymbol{\lambda}, \alpha)) + R_n(\boldsymbol{\psi}, \alpha), \quad (56)$$

where

$$\begin{aligned} B_n(\mathbf{t}_{\eta, \lambda}) &= 2\mathbf{t}_{\eta, \lambda}^\top \mathbf{S}_{\eta n} - \mathbf{t}_{\eta, \lambda}^\top \mathcal{I}_\eta \mathbf{t}_{\eta, \lambda}, \\ C_n(\mathbf{t}_\lambda) &= 2\mathbf{t}_\lambda^\top \mathbf{S}_{\lambda, \eta n} - \mathbf{t}_\lambda^\top \mathcal{I}_{\lambda, \eta} \mathbf{t}_\lambda \\ &= \mathbf{G}_{\lambda, \eta n}^\top \mathcal{I}_{\lambda, \eta} \mathbf{G}_{\lambda, \eta n} - (\mathbf{t}_\lambda - \mathbf{G}_{\lambda, \eta n})^\top \mathcal{I}_{\lambda, \eta} (\mathbf{t}_\lambda - \mathbf{G}_{\lambda, \eta n}), \end{aligned} \quad (57)$$

with $\mathbf{G}_{\lambda, \eta n} \xrightarrow{d} \mathbf{G}_{\lambda, \eta} = (\mathcal{I}_{\lambda, \eta})^{-1} \mathbf{S}_{\lambda, \eta}$ and $\mathbf{S}_{\lambda, \eta n} \xrightarrow{d} \mathbf{S}_{\lambda, \eta} \sim N(\mathbf{0}, \mathcal{I}_{\lambda, \eta})$. Also, $R_n(\hat{\psi}_\alpha, \alpha) = o_p(1)$ holds from Proposition 3(a) and $\sqrt{n}t(\hat{\psi}_\alpha, \alpha) = O_p(1)$.

Because $\Delta_{(\gamma, \theta)} f(x; \hat{\gamma}_0^*, \hat{\theta}_0^*)$ is identical to $\Delta_\eta f(x; \psi^*, \alpha)$, a standard analysis gives $2[L_{0, n}(\hat{\gamma}_0, \hat{\theta}_0) - L_{0, n}(\gamma_0^*, \theta_0^*)] = \max_{t_\eta} B_n(\sqrt{n}t_\eta) + o_p(1)$. Note that the possible values of both $\sqrt{n}t_\eta$ and $\sqrt{n}t_{\eta, \lambda}$ approaches \mathbb{R}^{q+2} . Therefore, in view of (56)-(57), we can write equation (55) as

$$LR_n = C_n(\sqrt{n}t_\lambda(\hat{\lambda}, \alpha)) + o_p(1), \quad (58)$$

where $\hat{\lambda}$ is defined in (14).

The asymptotic distribution of LR_n follows from applying Theorem 3(c) of (Andrews, 1999, p.1362) to (56) and (58). First, Assumption 2 of Andrews (1999) holds because Assumption 2* of Andrews (1999) hold because of Proposition 2(a). Second, Assumption 3 of Andrews (1999) holds with $B_T = n^{1/2}$ and $T = n$ because $\mathbf{S}_{\lambda, \eta n} \xrightarrow{d} \mathbf{S}_{\lambda, \eta} \sim N(\mathbf{0}, \mathcal{I}_{\lambda, \eta})$ and $\mathcal{I}_{\lambda, \eta}$ is non-singular. Assumption 4 of Andrews (1999) holds from part (a). Assumption 5 of Andrews (1999) follows from Assumption 5* and Lemma 3 of Andrews (1999) with $b_T = n^{1/2}$ because $\alpha(1 - \alpha)v(\Theta_\lambda)$ is locally equal to Λ_λ . Therefore, it follows from Theorem 3(c) of Andrews (1999) that $C_n(\sqrt{n}t_\lambda(\hat{\lambda}, \alpha)) \xrightarrow{d} (\hat{t}_\lambda)^\top \mathcal{I}_{\lambda, \eta} \hat{t}_\lambda$, where \hat{t}_λ is defined by (16). □

Proof of Proposition 5. Under $H_{2,0}$, we obtain the $\vartheta_{M_0+1} \in \Upsilon_{2h}^*$,

$$\begin{aligned} & \mathbb{E}[\{\nabla_{\alpha_h} \log f_{M_0+1}(\mathbf{W}_i, \vartheta_{M_0+1})\}^2] \\ &= \int \frac{\{f(\mathbf{w}; \boldsymbol{\mu}_h, \boldsymbol{\sigma}_h) - f(\mathbf{w}; \boldsymbol{\mu}_{M_0}^*, \boldsymbol{\sigma}_{M_0}^*)\}^2}{\sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)} d\mathbf{w} \\ &= \int \frac{\{f(\mathbf{w}; \boldsymbol{\mu}_h, \boldsymbol{\sigma}_h)\}^2}{\sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)} d\mathbf{w} + \int \frac{\{f(\mathbf{w}; \boldsymbol{\mu}_{M_0}^*, \boldsymbol{\sigma}_{M_0}^*)\}^2}{\sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)} d\mathbf{w} - 2 \int \frac{f(\mathbf{w}; \boldsymbol{\mu}_h, \boldsymbol{\sigma}_h) f(\mathbf{w}; \boldsymbol{\mu}_{M_0}^*, \boldsymbol{\sigma}_{M_0}^*)}{\sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)} d\mathbf{w}. \end{aligned} \quad (59)$$

The latter two terms on the right-hand side of (59) are bounded because $f(\mathbf{w}; \boldsymbol{\mu}_{M_0}^*, \boldsymbol{\sigma}_{M_0}^*) / \sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*) \leq (1/\alpha_{M_0}^*)$ for any \mathbf{w} and $f(\mathbf{w}; \mu, \sigma)$ integrates to one. Therefore, the left-hand side of (59) goes to infinity if and only if the first term on the right-hand side of (59) goes to infinity.

Because $\max_j \alpha_j \leq \sum_j^{M_0} \alpha_j \leq M_0 \max_j \alpha_j$, we obtain

$$\frac{1}{M_0 \max_j \{\alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)\}} \frac{\{f(\mathbf{w}; \boldsymbol{\mu}_h, \boldsymbol{\sigma}_h)\}^2}{\sum_{j=1}^{M_0} \alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)} \leq \frac{\{f(\mathbf{w}; \boldsymbol{\mu}_h, \boldsymbol{\sigma}_h)\}^2}{\max_j \{\alpha_j^* f(\mathbf{w}; \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^*)\}}.$$

Without loss of generality, we assume that $\sigma_{M_0}^* = \max\{\sigma_1^*, \dots, \sigma_{M_0}^2\}$ and the maximum is unique. Then there exists $M \in (0, \infty)$, such that $\max_j \{\alpha_j^* f(\mathbf{w}, \boldsymbol{\mu}_j^*, \boldsymbol{\sigma}_j^2)\} = \alpha_{M_0}^* f(\mathbf{w}, \boldsymbol{\mu}_{M_0}^*, \boldsymbol{\sigma}_{M_0}^2)$

when $|y_t| > M \forall t = 1, \dots, T$. Note that

$$\begin{aligned} \frac{\{f(\mathbf{w}; \boldsymbol{\mu}_h, \boldsymbol{\sigma}_h)\}^2}{f(\mathbf{w}; \boldsymbol{\mu}_{M_0}^*, \boldsymbol{\sigma}_{M_0}^*)} &= \prod_{t=1}^T \frac{(\sigma_{M_0}^*)^2}{(2\pi)^{1/2} \sigma_h^2} \exp \left\{ -\frac{1}{\sigma_h^2} (y_t - \mu_h)^2 + \frac{1}{2(\sigma_{M_0}^*)^2} (y_t - \mu_{M_0}^*)^2 \right\} \\ &= \frac{(\sigma_{M_0}^*)^{2T}}{(2\pi)^{1/2} \sigma_h^{2T}} \exp \left\{ -\frac{1}{\sigma_h^2} \sum_{t=1}^T (y_t - \mu_h)^2 + \frac{1}{2(\sigma_{M_0}^*)^2} \sum_{t=1}^T (y_t - \mu_{M_0}^*)^2 \right\}. \end{aligned} \quad (60)$$

The stated result follows because the integral of this over $|y| \geq M$ is finite if $\sigma_h^2/\sigma_{M_0}^{2*} < 2$ and infinite if $\sigma_h^2/\sigma_{M_0}^{2*} > 2$. When $\sigma_h^2/\sigma_{M_0}^{2*} = 2$, it is finite if $\mu_h = \mu_{M_0}^*$ and infinite if $\mu_h \neq \mu_{M_0}^*$. \square

Proof of Proposition 6. Our panel data model can be viewed as a special case of the T -dimensional multivariate normal mixture models, where the variance-covariance matrix for each component is given by a $T \times T$ diagonal matrix, $\Sigma_j := \text{diag}(\sigma_j^2, \dots, \sigma_j^2)$. Chen and Tan (2009) provides the consistency proof for the penalized maximum likelihood estimator for a multivariate normal mixture under their conditions C1-C3 for the penalty function but Alexandrovich (2014) identifies a soft spot in the proof of Chen and Tan (2009) and provides an alternative consistency proof by strengthening the condition C3 of Chen and Tan (2009). Their $p_n(G)$ and $\tilde{p}_n(G)$ correspond to our $\tilde{p}_n(\boldsymbol{\vartheta}_2)$ and $p_n(\boldsymbol{\vartheta}_2)$, respectively; consequently, the condition C1 and C2 in Chen and Tan (2009) and a version of condition C3 strengthened by Alexandrovich (2014) can be stated in our notation as follows:

C1. The penalty function is written as $\tilde{p}_n(\boldsymbol{\vartheta}_M) = \sum_{j=1}^M p_n(\sigma_j^2)$

C2. For any fixed $\boldsymbol{\vartheta}_M$ with $\sigma_j^2 > 0$ for $j = 1, 2, \dots, M$, we have $\tilde{p}_n(\boldsymbol{\vartheta}_M) = o(n)$ and $\sup_{\boldsymbol{\vartheta}_M \in \Theta_M} \max\{0, \tilde{p}_n(\boldsymbol{\vartheta}_M)\} = o(n)$. In addition, $\tilde{p}_n(\boldsymbol{\vartheta}_M)$ is differentiable with respect to $\boldsymbol{\vartheta}_M$ and as $n \rightarrow \infty$, $\nabla_{\boldsymbol{\vartheta}_M} \tilde{p}_n(\boldsymbol{\vartheta}_M) = o(n^{1/2})$ at any fixed $\boldsymbol{\vartheta}$ such that $\sigma_j^2 > 0$ for $j = 1, 2, \dots, M$.

A version of C3 by Alexandrovich (2014). For large enough n , $p_n(\sigma_j^2) \leq (\frac{3}{4}\sqrt{n \log \log n}) \log(\sigma_j^2)$, when $\sigma_j^2 < cn^{-2}$ for some $c > 0$.

The consistency of the PMLE, $\hat{\boldsymbol{\vartheta}}_{M_0}$ and $\hat{\boldsymbol{\vartheta}}_{M_0+1}$, follows from Theorems 1 and 3 of Chen and Tan (2009) and Corollary 3 of Alexandrovich (2014) if we can show that the above three conditions hold for our penalty function (3). Given (3), C1 trivially holds. Under Assumption 1(b), C2 also holds because $a_n = O(n^{1/4-\zeta})$ with $\zeta > 0$ implies $a_n = o(n)$ or $o(n^{1/2})$, and $\Delta_{\sigma_j^2} \tilde{p}_n(\boldsymbol{\vartheta}) = -a_n(-\sigma_0^2/(\sigma_j^2) + 1/\sigma_j^2) = o(n^{1/4})$ if $\sigma_j^2 > 0$. For C3, suppose that $\sigma_j^2 < n^{-2}$. Then, because $a_n = o(n^{1/4})$ and $a_n > 0$, $p_n(\sigma_j^2) = -a_n(\sigma_j^{-2}\sigma_0^2 + 2\log(\sigma_j/\sigma_0) - 1) < -c_n(n^{9/4}\sigma_0^2 - 2n^{5/4}\log(n\sigma_0) - n^{1/4}) < (\frac{3}{4}\sqrt{n \log \log n}) 2\log(n)$ when n is large enough, where c_n is a sequence of positive numbers that are bounded. Therefore, $\tilde{p}_n(\boldsymbol{\vartheta}_M)$ satisfies the above three conditions, and the stated result follows from Theorems 1 and 3 of Chen and Tan (2009) and Corollary 3 of Alexandrovich (2014). \square

Proof of Proposition 7. For $h = 1, \dots, M_0$, let $\mathcal{N}_h^* \subset \Theta_{\boldsymbol{\vartheta}_{M_0+1}}(\epsilon)$ be a sufficiently small closed neighbourhood of Υ_{1h}^* such that $\alpha_h, \alpha_{h+1} > 0$ hold and $\Upsilon_{1k}^* \not\subset \mathcal{N}_h^*$ if $k \neq h$. Consider the following one-to-one reparameterization from the $(M_0 + 1)$ -component model parameter $\boldsymbol{\vartheta}_{M_0+1} = (\alpha_1, \dots, \alpha_{M_0}, \boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_h^\top, \boldsymbol{\theta}_{h+1}^\top, \dots, \boldsymbol{\theta}_{M_0+1}^\top, \boldsymbol{\gamma}^\top)^\top$. Similarly to (5), the one-to-one reparameterization for testing the null hypothesis $H_{0,1h}$ is given by

$$\begin{pmatrix} \boldsymbol{\lambda}_h \\ \boldsymbol{\nu}_h \end{pmatrix} := \begin{pmatrix} \boldsymbol{\theta}_h - \boldsymbol{\theta}_{h+1} \\ \tau \boldsymbol{\theta}_h + (1 - \tau) \boldsymbol{\theta}_{h+1} \end{pmatrix} \text{ so that } \begin{pmatrix} \boldsymbol{\theta}_h \\ \boldsymbol{\theta}_{h+1} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\nu} + (1 - \tau) \boldsymbol{\lambda} \\ \boldsymbol{\nu} - \tau \boldsymbol{\lambda} \end{pmatrix}$$

and reparameterize α_j for $j = 1, 2, \dots, M_0$ as

$$\begin{aligned} (\pi_1, \dots, \pi_{h-1}, \pi_h, \pi_{h+1}, \dots, \pi_{M_0-1}) &= (\alpha_1, \dots, \alpha_{h-1}, (\alpha_h + \alpha_{h+1}), \alpha_{h+2}, \dots, \alpha_{M_0}) \\ \tau &= \alpha_h / (\alpha_h + \alpha_{h+1}) \end{aligned}$$

so that $\pi_h = \alpha_h + \alpha_{h+1}$ and $\pi_{M_0} = 1 - \sum_{j=1}^{M_0-1} \pi_j$.

Collect the reparameterized parameters except τ as

$$\boldsymbol{\psi}_{h,\tau} = (\boldsymbol{\eta}^\top, \boldsymbol{\lambda}_h^\top)^\top \quad \text{with} \quad \boldsymbol{\eta} = (\pi_1, \dots, \pi_{M_0-1}, \boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_{h-1}^\top, \boldsymbol{\nu}_h^\top, \boldsymbol{\theta}_{h+2}^\top, \dots, \boldsymbol{\theta}_{M_0+1}^\top, \boldsymbol{\gamma}^\top)^\top.$$

In the reparameterized model, the null restriction $\boldsymbol{\theta}_h = \boldsymbol{\theta}_{h+1}$ implied by $H_{0,1h}$ holds if and only if $\boldsymbol{\lambda}_h = 0$. Under $H_{0,1h}$, we have $\boldsymbol{\lambda}_h^* = 0$ and $\boldsymbol{\eta}^* = (\alpha_1^*, \dots, \alpha_{M_0-1}^*, (\boldsymbol{\theta}_1^*)^\top, \dots, (\boldsymbol{\theta}_{M_0}^*)^\top, (\boldsymbol{\gamma}^*)^\top)^\top$. Define the log-likelihood under the reparameterized parameters as

$$f_{M_0+1}^h(\mathbf{w}; \boldsymbol{\psi}_{h,\tau}, \tau) = \pi_h g^h(\mathbf{w}, \boldsymbol{\psi}_{h,\tau}, \tau) + \sum_{j=1}^{h-1} \pi_j f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\theta}_j) + \sum_{j=h}^{M_0} \pi_{j+1} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\theta}_{j+1}),$$

where $g^h(\mathbf{w}, \boldsymbol{\psi}_{h,\tau}, \tau)$ is defined similarly to (6) as

$$g^h(\mathbf{w}, \boldsymbol{\psi}_{h,\tau}, \tau) = \tau f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu}_h + (1 - \tau) \boldsymbol{\lambda}_h) + (1 - \tau) f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu}_h - \tau \boldsymbol{\lambda}_h). \quad (61)$$

Define the local penalized MLE of $\boldsymbol{\psi}_{h,\tau}$ by

$$\hat{\boldsymbol{\psi}}_{h,\tau} := \arg \max_{\boldsymbol{\psi}_{h,\tau} \in \mathcal{N}_h^*} L_n^h(\boldsymbol{\psi}_{h,\tau}, \tau) + \sum_{j=1}^{M_0} p_n(\sigma_j^2(\boldsymbol{\psi}_{h,\tau}, \tau)), \quad (62)$$

where $L_n^h(\boldsymbol{\psi}_{h,\tau}, \tau) := \sum_{i=1}^N \log g^h(\mathbf{W}_i; \boldsymbol{\psi}_{h,\tau}, \tau)$ and $\sigma_j^2(\boldsymbol{\psi}_{h,\tau}, \tau)$ is the value of σ_j^2 implied by the value of $\boldsymbol{\psi}_{h,\tau}$ and τ . Because $\boldsymbol{\psi}_{h,\tau}^*$ is the only parameter value in \mathcal{N}_h^* that generates the true density, $\hat{\boldsymbol{\psi}}_{h,\tau} - \boldsymbol{\psi}_{h,\tau}^* = o_p(1)$ follows Proposition 4.

For $\epsilon \in (0, 1/2)$, define the LRTS for testing $H_{0,1h}$ as $LR_n^{M_0,h} := \max_{\tau \in [\epsilon, 1-\epsilon]} 2(L_n^h(\hat{\boldsymbol{\psi}}_{h,\tau}, \tau) -$

$L_{0,n}(\hat{\vartheta}_{M_0})$) Because $\hat{\sigma}_j^2 - \sigma_{0,j}^2 = O_p(n^{-1/4})$ under the null hypothesis (c.f., Proposition 4(a)), we have $\tilde{p}_n(\vartheta_{M_0+1}) = o_p(1)$ by Assumption 3(c), and $PLR_n^{M_0,h} - LR_n^{M_0,h} = o_p(1)$ follows for $h = 1, \dots, M_0$.

Then, in view of (22), the stated result holds if

$$(LR_n^{M_0,1}, \dots, LR_n^{M_0,M_0})^\top \xrightarrow{d} (\hat{\mathbf{t}}_\lambda^1)^\top \mathcal{I}_{\eta,\lambda}(\hat{\mathbf{t}}_\lambda^1), \dots, (\hat{\mathbf{t}}_\lambda^{M_0})^\top \mathcal{I}_{\eta,\lambda}(\hat{\mathbf{t}}_\lambda^{M_0})^\top. \quad (63)$$

Observe that $L_n^h(\psi_{h,\tau}, \tau) - L_n^h(\psi_{h,\tau}^*, \tau)$ admits the same expansion as $L_n(\hat{\psi}, \alpha) - L_n(\psi^*, \alpha)$ in (13) and (56) by replacing $(\alpha, \mathbf{t}(\psi, \alpha), \mathbf{t}_\lambda(\lambda, \alpha), \mathbf{S}_n, \mathbf{G}_n, \mathcal{I}_n, R_n(\psi, \alpha))$ with $(\tau, \mathbf{t}^h(\psi^h, \tau), \mathbf{t}_\lambda^h(\lambda^h, \tau), \mathbf{S}_n^h, \mathbf{G}_n^h, \mathcal{I}_n^h, R_n^h(\psi^h, \tau))$, where $(\mathbf{S}_n^h, \mathcal{I}_n^h)$ is defined similarly to $(\mathbf{S}_n, \mathcal{I}_n)$ but replacing $(s_\eta, s_{\lambda\lambda})$ with $(\tilde{s}_\eta, s_{\lambda\lambda}^h)$ while $\mathbf{G}_n^h := (\mathcal{I}_n^h)^{-1} \mathbf{S}_n^h$. Applying the proof of Proposition 3, we have $\mathbf{S}_n^h \xrightarrow{d} \mathbf{S}^h \sim N(\mathbf{0}, \mathcal{I}^h)$ and $\mathcal{I}_n^h \xrightarrow{p} \mathcal{I}^h$. Then, (62) follows from the proof of Propositions 3 and 4 for each local penalized MLE by replacing $(\mathbf{G}_n, \hat{\mathbf{t}}_\lambda, \mathcal{I}_{\lambda,\eta})$ with $(\mathbf{G}_n^h, \hat{\mathbf{t}}_\lambda^h, \mathcal{I}_{\lambda,\eta}^h)$, and collecting the results while noting that $(\mathbf{S}_n^1, \dots, \mathbf{S}_n^{M_0}) \xrightarrow{d} (\mathbf{S}^1, \dots, \mathbf{S}^{M_0})$. \square

Proof of Proposition 8. The proof is similar to that of Proposition 7 in Kasahara and Shimotsu (2015). Let ω_n^h denote the sample counterpart of $(\hat{\mathbf{t}}_\lambda^h)^\top \mathcal{I}_{\lambda,\eta}^h \hat{\mathbf{t}}_\lambda^h$ in Proposition 7 such that the LRTS satisfies $2\{L_n^h(\hat{\psi}_{h,\tau}, \tau) - L_{0,n}(\hat{\vartheta}_{M_0})\} = \omega_n^h + o_p(1)$, where $\hat{\psi}_{h,\tau}$ is the local penalized MLE as defined (62) and ω_n^h is defined similarly to $C_n(\sqrt{n}\mathbf{t}_\lambda(\hat{\lambda}, \alpha))$ in (58) but replacing $(\mathbf{t}_\lambda(\hat{\lambda}, \alpha), \mathbf{S}_n, \mathbf{G}_n, \mathcal{I}_{\lambda,\eta})$ with $(\mathbf{t}_\lambda^h(\hat{\lambda}^h, \tau), \mathbf{S}_n, \mathbf{G}_n^h, \mathcal{I}_{\lambda,\eta}^h)$ defined in the proof of Proposition 7.

First, we show that $M_n^{h(1)}(\tau_0) = 2\{PL_n^h(\vartheta_{M_0+1}^{h(1)}(\tau_0), \tau_0) - L_{0,n}(\hat{\vartheta}_{M_0})\} = \omega_n^h + p(\tau_0) + o_p(1)$. Define $\vartheta_{M_0+1}^{h*}(\tau_0)$ by the value of ϑ_{M_0+1} in $\Theta_{\vartheta_{M_0+1}}^h(\tau_0) := \{\vartheta \in \Psi_h^* : \alpha_h/(\alpha_h + \alpha_{h+1}) = \tau_0\}$. Because $\vartheta_{M_0+1}^{h*}(\tau_0)$ is the only value of ϑ_{M_0+1} that yields the true density if $\vartheta_{M_0+1} \in \Psi_h^*$ in (20) and $\alpha_h/(\alpha_h + \alpha_{h+1}) = \tau_0$, $\vartheta_{M_0+1}^{h(1)}(\tau_0)$ equals a reparameterized penalized local MLE in the neighborhood of $\vartheta_{M_0+1}^{h*}(\tau_0)$, and $\vartheta_{M_0+1}^{h(1)}(\tau_0) - \vartheta_{M_0+1}^{h*}(\tau_0) = o_p(1)$ holds in view of Proposition 6. Furthermore, by the consistency of $\sigma_j^{h(1)}$ and $a_n = O(1)$, we have $\tilde{p}(\vartheta_{M_0+1}^{h(1)}(\tau_0)) \xrightarrow{p} 0$. Therefore, $M_n^{h(1)}(\tau_0) = \omega_n^h + p(\tau_0) + o_p(1)$ follows from repeating the proof of Proposition 7. Finally, $EM_n^{(1)} \xrightarrow{d} \max_{h=1}^{M_0} \{\omega_n^h\}$ holds because $\{0.5\} \in \mathcal{T}$ and $p(0.5) = 0$.

We proceed to show that $M_n^{h(K)}(\tau_0) = \omega_n^h + p(\tau_0) + o_p(1)$ for any finite K . Because a generalized EM step never decreases likelihood (Dempster et al., 1977), we have

$$PL_n(\vartheta_{M_0+1}^{h(K)}(\tau_0), \tau^{h(K)}(\tau_0)) > PL_n(\vartheta_{M_0+1}^{h(1)}(\tau_0), \tau^{h(1)}(\tau_0)). \quad (64)$$

Therefore, it follows from Theorem 1 of Chen and Tan (2009), Lemma 4 in Appendix B, and induction that $\vartheta_{M_0+1}^{h(K)}(\tau_0) - \vartheta_{M_0+1}^{h*} = o_p(1)$ for any finite K . Let $\tilde{\vartheta}_{M_0+1}^h$ be the maximizer of $PL_{M_0+1}(\vartheta_{M_0+1}, \tau^{h(K)}(\tau_0))$ under the constraint of $\alpha_h/(\alpha_h + \alpha_{h+1}) = \tau^{h(K)}(\tau_0)$ in an arbitrary small neighbourhood of $\vartheta_{M_0+1}^{h*}(\tau^{h(K)})$. Then, $2\{PL_n^h(\tilde{\vartheta}_{M_0+1}^h, \tau^{h(K)}(\tau_0)) - L_{0,n}(\hat{\vartheta}_{M_0})\} = \omega_n^h + p(\tau_0) + o_p(1)$ holds from the definition of $\tilde{\vartheta}_{M_0+1}^h$ and $\tilde{p}(\tilde{\vartheta}_{M_0+1}^h) \xrightarrow{p} 0$ by repeating the proof

of Proposition 7. It also follows from the consistency of $\boldsymbol{\vartheta}_{M_0+1}^{h(K)}(\tau_0)$ that $PL_n(\tilde{\boldsymbol{\vartheta}}_n^h, \tau^{h(K)}(\tau_0)) \geq PL_n(\boldsymbol{\vartheta}_{M_0+1}^{h(K)}(\tau_0), \tau^{h(K)}(\tau_0)) + o_p(1)$. Therefore, in view of (64), we have

$$PL_n(\tilde{\boldsymbol{\vartheta}}_n^h, \tau^{h(K)}(\tau_0)) \geq PL_n(\boldsymbol{\vartheta}_{M_0+1}^{h(K)}(\tau_0), \tau^{h(K)}(\tau_0)) + o_p(1) \geq PL_n(\boldsymbol{\vartheta}_{M_0+1}^{h(1)}(\tau_0), \tau^{h(1)}(\tau_0)). \quad (65)$$

Finally, because $2\{PL_n(\tilde{\boldsymbol{\vartheta}}_{M_0+1}^h, \tau^{h(K)}(\tau_0)) - L_{0,n}(\hat{\boldsymbol{\vartheta}}_{M_0})\} = \omega_n^h + p(\tau_0) + o_p(1)$ and $2\{PL_n(\boldsymbol{\vartheta}_{M_0+1}^{h(1)}(\tau_0), \tau^{h(1)}(\tau_0)) - L_{0,n}(\hat{\boldsymbol{\vartheta}}_{M_0})\} = \omega_n^h + p(\tau_0) + o_p(1)$, it follows from (65) that $M_n^{h(K)}(\tau_0) = 2\{PL_n(\boldsymbol{\vartheta}_{M_0+1}^{h(K)}(\tau_0)) - L_{0,n}(\hat{\boldsymbol{\vartheta}}_{M_0})\} = \omega_n^h + p(\tau_0) + o_p(1)$ holds for all h . The stated result then follows from the definition of $EM_n^{(K)}$ and $\{0.5\} \in \mathcal{T}$. \square

Proof of Proposition 9. Let $\boldsymbol{\psi}_n = ((\boldsymbol{\nu}^*)^\top, \boldsymbol{\lambda}_n^\top)^\top$ be the value of $\boldsymbol{\psi}$ under $H_{1,n} : \boldsymbol{\vartheta} = \boldsymbol{\vartheta}_{2,n}$ and let $\boldsymbol{h} = (\mathbf{0}^\top, \boldsymbol{h}_\lambda^\top)^\top$, where \boldsymbol{h}_λ is defined by (31). Let $\mathbb{P}_{\boldsymbol{\vartheta}}$ be the probability measure on $\{\mathbf{W}_i\}_{i=1}^n$ under $\boldsymbol{\vartheta}$. Denote the log-likelihood ratio of $\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}$ to $\mathbb{P}_{\boldsymbol{\vartheta}_2^*}$ by $\log\left(\frac{d\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}}{d\mathbb{P}_{\boldsymbol{\vartheta}_2^*}}\right) = L_n(\boldsymbol{\psi}_n, \alpha^*) - L_n(\boldsymbol{\psi}^*, \alpha^*)$. Then, it follows from (12) and Proposition 3 that

$$\log\frac{d\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}}{d\mathbb{P}_{\boldsymbol{\vartheta}_2^*}} = \boldsymbol{h}^\top \boldsymbol{S}_n - \boldsymbol{h}^\top \boldsymbol{\mathcal{I}} \boldsymbol{h} / 2 + o_p(1) \quad \text{under } \mathbb{P}_{\boldsymbol{\vartheta}_2^*}. \quad (66)$$

Furthermore, because $\boldsymbol{S}_n \xrightarrow{d} N(\mathbf{0}, \boldsymbol{\mathcal{I}})$ under $\mathbb{P}_{\boldsymbol{\vartheta}_2^*}$, $\frac{d\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}}{d\mathbb{P}_{\boldsymbol{\vartheta}_2^*}}$ converges in distribution under $\mathbb{P}_{\boldsymbol{\vartheta}_2^*}$ to $\exp(N(\boldsymbol{\mu}, \boldsymbol{\sigma}^2))$ with $\boldsymbol{\mu} = -(1/2)\boldsymbol{h}^\top \boldsymbol{\mathcal{I}} \boldsymbol{h}$ and $\boldsymbol{\sigma}^2 = \boldsymbol{h}^\top \boldsymbol{\mathcal{I}} \boldsymbol{h}$ so that $E(\exp(N(\boldsymbol{\mu}, \boldsymbol{\sigma}^2))) = 1$. Consequently, $\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}$ is mutually contiguous with respect to $\mathbb{P}_{\boldsymbol{\vartheta}_2^*}$ from Le Cam's First Lemma (see, e.g., Corollary 12.3.1 of Lehmann and Romano, 2005) and, in view of (66), we have

$$\left(\log \frac{d\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}}{d\mathbb{P}_{\boldsymbol{\vartheta}_2^*}} \right) \xrightarrow{d} N \left(\begin{pmatrix} \mathbf{0} \\ -(1/2)\boldsymbol{h}^\top \boldsymbol{\mathcal{I}} \boldsymbol{h} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\mathcal{I}} & \boldsymbol{\mathcal{I}} \boldsymbol{h} \\ \boldsymbol{h}^\top \boldsymbol{\mathcal{I}} & \boldsymbol{h}^\top \boldsymbol{\mathcal{I}} \boldsymbol{h} \end{pmatrix} \right) \quad \text{under } \mathbb{P}_{\boldsymbol{\vartheta}_2^*}.$$

and

$$\boldsymbol{S}_n \xrightarrow{d} N(\boldsymbol{\mathcal{I}} \boldsymbol{h}, \boldsymbol{\mathcal{I}}) \quad \text{under } \mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}$$

from Le Cam's Third Lemma (see, e.g., 12.3.2 of Lehmann and Romano, 2005). Therefore, the proof of Proposition 4 goes through under $\mathbb{P}_{\boldsymbol{\vartheta}_{2,n}}$ if we replace $\boldsymbol{S}_{\lambda, \eta n} \xrightarrow{d} \boldsymbol{S}_{\lambda, \eta}$ with $\boldsymbol{S}_{\lambda, \eta n} \xrightarrow{d} \boldsymbol{S}_{\lambda, \eta} + (\boldsymbol{\mathcal{I}}_\lambda - \boldsymbol{\mathcal{I}}_{\eta\lambda} \boldsymbol{\mathcal{I}}_\eta^{-1} \boldsymbol{\mathcal{I}}_{\eta\lambda}) \boldsymbol{h}_\lambda = \boldsymbol{S}_{\lambda, \eta} + \boldsymbol{\mathcal{I}}_{\lambda, \eta} \boldsymbol{h}_\lambda$, and the stated result follows. \square

Proof of Proposition 10. We provide a proof for \hat{M}_{PLR} . The consistency proof for \hat{M}_{EM} is similar. We first prove that, when $M < M_0$, $\Pr(PLR_n(M) > \hat{c}_{1-q_n}^M) \rightarrow 1$ as $n \rightarrow \infty$. Let $\tilde{Q}_n^M(\boldsymbol{\vartheta}_M) := Q_n^M(\boldsymbol{\vartheta}_M) + n^{-1} \tilde{p}_n(\boldsymbol{\vartheta}_M)$ for $M \geq 2$. By Assumptions 3(c) and 4(b), $n^{-1} \tilde{p}_n(\boldsymbol{\vartheta}_M) \xrightarrow{P} 0$ uniformly over $\Theta_{\boldsymbol{\vartheta}_M}$. Then, it follows from Lemma 2.4 of Newey and McFadden (1994) that

$$\sup_{\boldsymbol{\vartheta}_M \in \Theta_{\boldsymbol{\vartheta}_M}} |\tilde{Q}_n^M(\boldsymbol{\vartheta}_M) - Q_n^M(\boldsymbol{\vartheta}_M)| = o_p(1), \quad (67)$$

and Assumption 4(a)(b) and the standard consistency proof (e.g., Theorem 2.1 of Newey and McFadden, 1994) give $\hat{\boldsymbol{\vartheta}}_M \xrightarrow{P} \boldsymbol{\vartheta}_M^*$ for $M < M_0$. Furthermore, by Assumption 3(c), $n^{-1/4} \nabla \tilde{p}_n(\boldsymbol{\vartheta}_M) = o_p(1)$ uniformly over $\Theta_{\boldsymbol{\vartheta}_M}$, and it follows from the argument in Theorem 3.2 of White (1982) that

$$\sqrt{n}(\hat{\boldsymbol{\vartheta}}_M - \boldsymbol{\vartheta}_M^*) \xrightarrow{P} N(0, A^M(\boldsymbol{\vartheta}_M^*)^{-1} B^M(\boldsymbol{\vartheta}_M^*) A^M(\boldsymbol{\vartheta}_M^*)^{-1}). \quad (68)$$

Then, from (67), (68), and the mean value expansion, we have $\tilde{Q}_n^M(\hat{\boldsymbol{\vartheta}}_M) - Q^M(\boldsymbol{\vartheta}_M^*) = O_p(n^{-1/2})$, and

$$\frac{PLR_n(M)}{n} := \tilde{Q}_n^{M+1}(\hat{\boldsymbol{\vartheta}}_{M+1}) - \tilde{Q}_n^M(\hat{\boldsymbol{\vartheta}}_M) = Q^{M+1}(\boldsymbol{\vartheta}_{M+1}^*) - Q^M(\boldsymbol{\vartheta}_M^*) + o_p(1).$$

Because $Q^{M+1}(\boldsymbol{\vartheta}_{M+1}^*) - Q^M(\boldsymbol{\vartheta}_M^*) > 0$ by Assumption 4(f), $PLR_n(M)/n \rightarrow \infty$ as $n \rightarrow \infty$. By Lemma 2, $-n^{-1} \ln q_n = o(1)$ and $\hat{c}_{1-q_n}^M - c_{1-q_n}^M = o_p(1)$ implies that $n^{-1} \hat{c}_{1-q_n}^M = o_p(1)$. Therefore, when $M < M_0$, we have $\Pr(PLR_n(M) > \hat{c}_{1-q_n}^M) = \Pr(PLR_n(M)/n > \hat{c}_{1-q_n}^M/n) \rightarrow 1$ as $n \rightarrow \infty$.

When $M = M_0$, because $PLR_n(M_0) = O_p(1)$ by Proposition 7 and $\hat{c}_{1-q_n}^{M_0} \rightarrow \infty$ by $q_n = o(1)$, $\Pr(PLR_n(M_0) > \hat{c}_{1-q_n}^{M_0}) \rightarrow 0$ as $n \rightarrow \infty$. \square

B Auxiliary results and their proofs

B.1 Lemmas

Lemma 1. For any $M < \infty$, $\Pr(-\log n + \ell(\mathbf{W}_{i^*}; \bar{Y}_{i^*}, s_{i^*}^2) < M) \rightarrow 0$ as $n \rightarrow \infty$.

Proof of Lemma 1. Because $\sum_{t=1}^T \frac{(Y_{it} - \mu)^2}{s_{i^*}^2} = T - 1$ when $i = i^*$, we have

$$\begin{aligned} -\log n + \ell(\mathbf{W}_{i^*}; \bar{Y}_{i^*}, s_{i^*}^2) &= -\log n - \frac{T}{2} \log s_{i^*}^2 - \frac{T}{2} \log(2\pi) - \frac{T-1}{2} \\ &= -\log \left(Cn(s_{i^*}^2)^{T/2} \right), \end{aligned} \quad (69)$$

for some positive constant C .

Therefore, to prove the stated result, it suffices to show that, for any $\epsilon > 0$, $\Pr(n(s_{i^*}^2)^{T/2} > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$. Given the property of the first-order statistic, the distribution of $s_{i^*}^2$ is given by $1 - [1 - F_{T-1}(s)]^n$, where $F_{T-1}(s)$ is the cumulative distribution function for chi-squared variables of degree $T - 1$. It follows that

$$\Pr(n(s_{i^*}^2)^{T/2} > \epsilon) = [1 - F_{T-1}((\epsilon/n)^{2/T})]^n.$$

When $T = 3$, $1 - F_{T-1}(s) = e^{-s/2}$ and $\Pr(n(s_{i^*}^2)^{T/2} > \epsilon) = e^{-Cn^{1/3}}$ for some positive constant C and, therefore, $\Pr(n(s_{i^*}^2)^{T/2} > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$, and the stated result follows.

For general $T \geq 2$, write

$$\left[1 - F_{T-1}((\epsilon/n)^{2/T})\right]^n = \left\{ \left[1 - F_{T-1}((\epsilon/n)^{2/T})\right]^{\frac{1}{F_{T-1}((\epsilon/n)^{2/T})}} \right\}^{n F_{T-1}((\epsilon/n)^{2/T})}. \quad (70)$$

Then, because $(1 - F)^{\frac{1}{F}} \rightarrow \frac{1}{e}$ when $F \rightarrow 0$, the stated result follows from (70) if we can show

$$\frac{F_{T-1}((\epsilon x)^{2/T})}{x} \rightarrow \infty \text{ as } x \rightarrow 0$$

for $x = 1/n$. By applying L'Hôpital's rule, we have

$$\lim_{x \rightarrow 0} \frac{F_{T-1}((\epsilon x)^{2/T})}{x} = \lim_{x \rightarrow 0} f_{T-1}((\epsilon x)^{2/T}) \epsilon^{2T} x^{2/T-1},$$

where f_k is the PDF of χ -square distribution with k degrees of freedom. Note that $f_{T-1}((\epsilon x)^{2/T}) = \frac{1}{2^{(T-1)/2} \Gamma((T-1)/2)} ((\epsilon x)^{2/T})^{(T-1)/2-1} e^{-((\epsilon x)^{2/T})/2}$, then

$$\lim_{x \rightarrow 0} f_{T-1}((\epsilon x)^{2/T}) x^{2/T-1} = \lim_{x \rightarrow 0} C_{T,\epsilon} e^{-((\epsilon x)^{2/T})/2} x^{-\frac{1}{T}} = \infty,$$

where $C_{T,\epsilon} = \frac{\epsilon^{(T-1)/T}}{2^{(T-1)/2} \Gamma((T-1)/2)}$, because $e^{-((\epsilon x)^{2/T})/2} \rightarrow 1$ and $x^{-\frac{1}{T}} \rightarrow \infty$ as $x \rightarrow 0$ for any finite $T \geq 2$. Therefore, $\lim_{x \rightarrow 0} \frac{F_{T-1}((\epsilon x)^{2/T})}{x} = \infty$ and the stated result for $T \geq 2$ follows from (70). \square

Lemma 2. Suppose that assumptions in Proposition 10 hold. If $-n^{-1} \ln q_n = o(1)$, then $n^{-1} c_{1-q_n}^M = o(1)$.

Proof. For brevity of notation, write $c_n = c_{1-q_n}^M$. By Theorem 2.1 of Foutz and Srivastava (1977), $PLR_n(M) \xrightarrow{d} \sum_{j=1}^K b_j \chi_j^2$ for $0 < b_j < \infty$ and K is finite, where $\chi_1^2, \dots, \chi_K^2$ are independent chi-square random variables with one degree of freedom. Then, we have

$$q_n = \Pr \left(\sum_{j=1}^K b_j \chi_j^2 \geq c_n \right) \leq \sum_{j=1}^K \Pr \left(\chi_j^2 \geq \frac{c_n}{b_j} \right) \leq \frac{K}{\sqrt{1-2t}} \exp \left(-t \frac{c_n}{b^*} \right) \quad \text{for } 0 < t < \frac{1}{2}$$

with $b^* = \arg \max \{b_1, \dots, b_K\}$, where the last inequality follows from a Chernoff bound: $\Pr \left(\chi_j^2 \geq \frac{c_n}{b^*} \right) \leq \frac{\mathbb{E}[\exp(t(\chi_j^2-1))]}{\exp(t(\frac{c_n}{b^*}-1))} = \frac{1}{\sqrt{1-2t}} \exp \left(-t \frac{c_n}{b^*} \right)$ for $0 < t < \frac{1}{2}$. Therefore, $-\frac{\ln q_n}{n} \geq -\frac{1}{n} \ln \left(\frac{K}{\sqrt{1-2t}} \right) + \frac{1}{2b^*} \frac{c_n}{n}$, and the stated result follows. \square

Lemma 3. Suppose that $g(\mathbf{w}; \boldsymbol{\psi}, \alpha)$ is defined as (6), where $\boldsymbol{\psi} = (\boldsymbol{\eta}^\top, \boldsymbol{\lambda}^\top)^\top$. Let $g^*, \nabla g^*, \nabla \log g^*$ denote $g(\mathbf{W}; \boldsymbol{\psi}, \alpha), \nabla g(\mathbf{W}; \boldsymbol{\psi}, \alpha)$, and $\nabla \log g(\mathbf{W}; \boldsymbol{\psi}, \alpha)$ evaluated at $(\boldsymbol{\psi}, \alpha)$, respectively. Let ∇f^* denote $\nabla f(\mathbf{W}; \boldsymbol{\gamma}^*, \boldsymbol{\theta}^*)$. The following statements hold.

- (a) For $l = 0, 1, \dots$, $\nabla_{(\boldsymbol{\lambda} \otimes \boldsymbol{\eta}^{\otimes l})^\top} g^* = 0$;
- (b) $\nabla_{(\boldsymbol{\lambda}^{\otimes 2})^\top} g^* = \alpha(1 - \alpha) \nabla_{(\boldsymbol{\theta}^{\otimes 2})^\top} f^*$;
- (c) $\mathbb{E}[\nabla_{\lambda_i \lambda_j} \log g^*] = 0$, $\mathbb{E}[\nabla_{\lambda_i \lambda_j \lambda_k} \log g^*] = 0$, and $\mathbb{E}[\nabla_{\eta \lambda_i \lambda_j} \log g^*] = -\mathbb{E}[\nabla_{\eta} \log g^* \nabla_{\lambda_i \lambda_j} \log g^*]$;
- (d) $\mathbb{E}[\nabla_{\lambda_i \lambda_j \lambda_k \lambda_\ell} \log g^*] = -\mathbb{E}[\nabla_{\lambda_i \lambda_j} \log g^* \nabla_{\lambda_k \lambda_\ell} \log g^* + \nabla_{\lambda_i \lambda_k} \log g^* \nabla_{\lambda_j \lambda_\ell} \log g^* + \nabla_{\lambda_i \lambda_\ell} \log g^* \nabla_{\lambda_j \lambda_k} \log g^*]$.

Proof of Lemma 3. Recall that

$$g(\mathbf{w}; \boldsymbol{\psi}, \alpha) = \alpha f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} + (1 - \alpha)\boldsymbol{\lambda}) + (1 - \alpha)f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda}).$$

First we show that for $l = 0$ holds for (a), $\nabla_{\boldsymbol{\lambda}} g^* = \alpha(1 - \alpha) \nabla_{\boldsymbol{\theta}} f^* - \alpha(1 - \alpha) \nabla_{\boldsymbol{\theta}} f^* = 0$. For $l > 0$, by Fubini's Theorem, we have

$$\begin{aligned} \nabla_{(\boldsymbol{\lambda}^{\otimes 2})^\top} g &= \nabla_{\boldsymbol{\lambda}} \left(\alpha \nabla_{(\boldsymbol{\gamma}, \boldsymbol{\theta})^{\otimes l}} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} + (1 - \alpha)\boldsymbol{\lambda}) + (1 - \alpha) \nabla_{(\boldsymbol{\gamma}, \boldsymbol{\theta})^{\otimes l}} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda}) \Big|_{\boldsymbol{\nu}=\boldsymbol{\theta}^*, \boldsymbol{\lambda}=\mathbf{0}} \right) \\ &= \left(\alpha(1 - \alpha) \nabla_{(\boldsymbol{\gamma}^{\otimes l}, \boldsymbol{\theta}^{\otimes l+1})} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} + (1 - \alpha)\boldsymbol{\lambda}) - \alpha(1 - \alpha) \nabla_{(\boldsymbol{\gamma}^{\otimes l}, \boldsymbol{\theta}^{\otimes l+1})} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda}) \Big|_{\boldsymbol{\nu}=\boldsymbol{\theta}^*, \boldsymbol{\lambda}=\mathbf{0}} \right) \\ &= 0. \end{aligned}$$

To show part (b), note that

$$\begin{aligned} \nabla_{(\boldsymbol{\lambda}^{\otimes 2})^\top} g &= \nabla_{\boldsymbol{\lambda}} \left(\alpha(1 - \alpha) \nabla_{\boldsymbol{\lambda}^\top} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha(1 - \alpha)\boldsymbol{\lambda}) + (1 - \alpha) \nabla_{\boldsymbol{\lambda}^\top} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda}) \right) \\ &= \alpha(1 - \alpha)^2 \nabla_{(\boldsymbol{\lambda}^{\otimes 2})^\top} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} + \alpha^2(1 - \alpha)\boldsymbol{\lambda}) + (1 - \alpha) \nabla_{(\boldsymbol{\lambda}^{\otimes 2})^\top} f(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\nu} - \alpha\boldsymbol{\lambda}) \Big|_{\boldsymbol{\nu}=\boldsymbol{\theta}^*, \boldsymbol{\lambda}=\mathbf{0}} \\ &= \nabla_{(\boldsymbol{\lambda}^{\otimes 2})^\top} f^*. \end{aligned}$$

For parts (c) and (d), observe that $\int \nabla_{\lambda_i} \log g(\mathbf{w}; \boldsymbol{\psi}, \alpha) g(\mathbf{w}; \boldsymbol{\psi}, \alpha) dx = 0$ holds for any $\boldsymbol{\psi}$ in the interior of $\Theta_{\boldsymbol{\psi}}$, and differentiating this equation w.r.t. λ_j gives

$$\int \{ \nabla_{\lambda_i \lambda_j} \log g(\mathbf{w}; \boldsymbol{\psi}, \alpha) + \nabla_{\lambda_i} \log g(\mathbf{w}; \boldsymbol{\psi}, \alpha) \nabla_{\lambda_j} \log g(\mathbf{w}; \boldsymbol{\psi}, \alpha) \} g(\mathbf{w}; \boldsymbol{\psi}, \alpha) dx = 0. \quad (71)$$

Evaluating (71) at $\boldsymbol{\psi} = \boldsymbol{\psi}^*$ in conjunction with part (a) gives the first equation in part (c). Differentiating (71) w.r.t. λ_k or η and evaluating at $\boldsymbol{\psi} = \boldsymbol{\psi}^*$ give the latter two equations in part (c). Part (d) follows from differentiating (71) w.r.t. λ_k and λ_ℓ and evaluating at $\boldsymbol{\psi} = \boldsymbol{\psi}^*$ in conjunction with parts (a)(c). □

Lemma 4. *Suppose that the assumptions of Proposition 8 hold. If $\boldsymbol{\vartheta}_{M_0+1}^{h(K)}(\tau_0) - \boldsymbol{\vartheta}_{M_0+1}^{h^*}(\tau_0) = o_p(1)$ and $\tau^{(K)} - \tau_0 = o_p(1)$, then (a) $\alpha_m^{(K+1)} / [\alpha_h^{(K+1)} + \alpha_{h+1}^{(K+1)}] - \tau_0 = o_p(1)$ and (b) $\tau^{(K+1)} - \tau_0 = o_p(1)$.*

Proof. The proof is similar to the proof of Lemma 3 of Chen and Li (2009) and Lemma 10 in Appendix D of Kasahara and Shimotsu (2019). We suppress (τ_0) from $\boldsymbol{\vartheta}_{M_0+1}^{h(K)}(\tau_0)$ and $\boldsymbol{\vartheta}_{M_0+1}^{h*}(\tau_0)$. We suppress \mathbf{Z} for brevity. Let $f_i(\boldsymbol{\gamma}, \boldsymbol{\theta}_j)$ and $f_i(\boldsymbol{\vartheta}_{M_0+1})$ denote $f(\mathbf{W}_i; \boldsymbol{\gamma}, \boldsymbol{\theta}_j)$ in (2) and $f_{M_0+1}(\mathbf{W}_i; \boldsymbol{\vartheta}_{M_0+1})$ in (18), respectively. Applying a Taylor expansion to $\alpha_h^{(K+1)} = n^{-1} \sum_{i=1}^n w_{ih}^{(K)}$ and using $\boldsymbol{\vartheta}_{M_0+1}^{h(K)} - \boldsymbol{\vartheta}_{M_0+1}^{h*} = o_p(1)$, we obtain

$$\begin{aligned}\alpha_m^{(K+1)} &= \frac{1}{n} \sum_{i=1}^n \frac{\tau^{(K)}(\alpha_h^{(K)} + \alpha_{h+1}^{(K)}) f_i(\boldsymbol{\gamma}^{(K)}, \boldsymbol{\theta}_h^{(K)})}{f_i(\boldsymbol{\vartheta}_{M_0+1}^{h(K)})} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\tau_0 \alpha_h^* f_i(\boldsymbol{\gamma}^*, \boldsymbol{\theta}_h^*)}{f_i(\boldsymbol{\vartheta}_{M_0+1}^{h*})} + o_p(1) = \tau_0 \alpha_h^* + o_p(1),\end{aligned}$$

where the last equality follows from $\mathbb{E}[f_i(\boldsymbol{\gamma}^*, \boldsymbol{\theta}_h^*)/f_i(\boldsymbol{\vartheta}_{M_0+1}^{h*})] = 1$ and the law of large numbers. A similar argument gives $\alpha_{h+1}^{(K+1)} = (1 - \tau_0) \alpha_{h+1}^* + o_p(1)$, and part (a) follows.

For part (b), define $H(\tau) := \sum_{i=1}^n w_{ih}^{(K)} \log(\tau) + \sum_{i=1}^n w_{i,h+1}^{(K)} \log(1 - \tau) = n \alpha_h^{(K+1)} \log(\tau) + n \alpha_{h+1}^{(K+1)} \log(1 - \tau)$, then $\tau^{(K+1)}$ maximizes $H(\tau) + p(\tau)$. $H(\tau)$ is maximized at $\tilde{\tau} := \alpha_h^{(K+1)} / (\alpha_h^{(K+1)} + \alpha_{h+1}^{(K+1)}) = (\tau_0 \alpha_h^* + o_p(1)) / (\tau_0 \alpha_h^* + (1 - \tau_0) \alpha_{h+1}^* + o_p(1)) = \tau_0 + o_p(1)$. Observe that, with $\tilde{\tau}$ between $\tau^{(K+1)}$ and $\tilde{\tau}$,

$$p(\tilde{\tau}) \leq p(\tau) - p(\tau^{(K+1)}) \leq H(\tau^{(K+1)}) - H(\tilde{\tau}) = H''(\tilde{\tau})(\tau^{(K+1)} - \tilde{\tau})^2, \quad (72)$$

where the first inequality follows from $p(\tau) \leq 0$, the second inequality holds because $\tau^{(K+1)}$ maximizes $H(\tau) + p(\tau)$, and the last equality follows from expanding $H(\tau^{(K+1)})$ twice around $\tilde{\tau}$ and noting that $H'(\tilde{\tau}) = 0$ because $\tilde{\tau}$ maximizes $H(\tau)$. Note that $H''(\tau) = -n \times \left\{ \frac{\alpha_h^{(K+1)}}{\tau^2} + \frac{\alpha_{h+1}^{(K+1)}}{(1-\tau)^2} \right\} < 0$ and $\inf_{\tau} H''(\tau) \geq -n(\alpha_h^{(K+1)} + \alpha_{h+1}^{(K+1)})$. Therefore, in view of $\tilde{\tau} - \tau_0 = o_p(1)$ and (72), we have $(\tau^{(K+1)} - \tilde{\tau})^2 \leq p(\tilde{\tau})/H''(\tilde{\tau}) = O_p(n^{-1})$, and part (b) holds. \square

B.2 Score function for testing $H_0 : m = 1$ against $H_A : m = 2$

$H^j(\cdot)$ is defined as the j -th order Hermite polynomial. $H^1(t) = t$, $H^2(t) = t^2 - 1$, $H^3(t) = t^3 - 3t$, and $H^4(t) = t^4 - 6t^2 + 3$. As shown in Kasahara and Shimotsu (2015) supplement material, the derivative of $\{\frac{1}{\sigma} \phi(\frac{t}{\sigma})\}$ is

$$\frac{\nabla_{\mu^m} \nabla_{(\sigma^2)^\ell} \{\frac{1}{\sigma} \phi(\frac{t}{\sigma})\}}{\{\frac{1}{\sigma} \phi(\frac{t}{\sigma})\}} = \left(\frac{1}{2}\right)^\ell \left(\frac{1}{\sigma}\right)^{m+2\ell} H^{m+2\ell} \left(\frac{t}{\sigma}\right).$$

Let

$$f^* = f(\mathbf{w}; \gamma^*, \theta^*), \nabla f^* = \nabla f(\mathbf{w}; \gamma^*, \theta^*), H_{i,t}^{j*} = \frac{1}{\sigma^{*j} j!} H^j \left(\frac{y_{it} - \mathbf{x}_{it}^\top \boldsymbol{\beta}^* - \mathbf{z}_{it}^\top \boldsymbol{\gamma}^* - \mu^*}{\sigma^*} \right), \quad (73)$$

then the first-order derivatives of the density functions are

$$\begin{aligned} \nabla_\mu f^* &= f^* \sum_{t=1}^T \frac{1}{\sigma} H_{i,t}^{1*}; \quad \nabla_{\sigma^2} f^* = f^* \sum_{t=1}^T \frac{1}{2} \frac{1}{\sigma^2} H_{i,t}^{2*}; \\ \nabla_\beta f^* &= f^* \sum_{t=1}^T \frac{1}{\sigma} H_{i,t}^{1*} \mathbf{x}_{it}; \quad \nabla_\gamma f^* = f^* \sum_{t=1}^T \frac{1}{\sigma} H_{i,t}^{1*} \mathbf{z}_{it}. \end{aligned}$$

The score function defined in (9) is then written in terms of the Hermite polynomials:

$$\mathbf{s}_\eta = \begin{pmatrix} s_\mu \\ s_\sigma \\ \mathbf{s}_\beta \\ s_\gamma \end{pmatrix} = \begin{pmatrix} \sum_{t=1}^T H_{i,t}^{1*} \\ \sum_{t=1}^T H_{i,t}^{2*} \\ \sum_{t=1}^T H_{i,t}^{1*} \mathbf{x}_{it} \\ \sum_{t=1}^T H_{i,t}^{1*} \mathbf{z}_{it} \end{pmatrix}, \quad \mathbf{s}_{\lambda\lambda} = \begin{pmatrix} s_{\lambda_\mu \lambda_\mu} \\ s_{\lambda_\mu \lambda_\sigma} \\ s_{\lambda_\sigma \lambda_\sigma} \\ s_{\lambda_\mu \lambda_\beta} \\ s_{\lambda_\sigma \lambda_\beta} \\ s_{\lambda_\beta \lambda_\beta} \end{pmatrix}, \quad (74)$$

where

$$\begin{aligned} \begin{pmatrix} s_{\lambda_{\mu\mu}} \\ s_{\lambda_{\mu\sigma}} \\ s_{\lambda_{\sigma\sigma}} \\ s_{\lambda_{\mu\beta}} \\ s_{\lambda_{\sigma\beta}} \end{pmatrix} &= \begin{pmatrix} \sum_{t=1}^T H_{i,t}^{2*} + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{1,i,t}^{1*} H_{i,s}^{1*} \\ 3 \sum_{t=1}^T H_{i,t}^{3*} + \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} H_{i,s}^{2*} \\ 3 \sum_{t=1}^T H_{i,t}^{4*} + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{2*} H_{i,t}^{2*} \\ 2 \sum_{t=1}^T H_{i,t}^{2*} \mathbf{x}_{it} + \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} H_{i,s}^{1*} \mathbf{x}_{it} \\ 3 \sum_{t=1}^T H_{i,t}^{3*} \mathbf{x}_{it} + 2 \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} H_{i,s}^{2*} \mathbf{x}_{it} \end{pmatrix}, \quad \text{and} \\ s_{\lambda_{\beta\beta}} &= \begin{pmatrix} \sum_{t=1}^T H_{i,t}^{2*} x_{it,1}^2 + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} x_{it,1} H_{i,s}^{1*} x_{is,1} \\ \vdots \\ \sum_{t=1}^T H_{i,t}^{2*} x_{it,q}^2 + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} x_{it,q} H_{i,s}^{1*} x_{is,q} \\ 2 \sum_{t=1}^T H_{i,t}^{2*} x_{it,1} x_{it,2} + \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} x_{it,1} H_{i,s}^{1*} x_{is,2} \\ \vdots \\ 2 \sum_{t=1}^T H_{i,t}^{2*} x_{it,1} x_{it,q} + \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} x_{it,1} H_{i,s}^{1*} x_{is,q} \\ 2 \sum_{t=1}^T H_{i,t}^{2*} x_{it,2} x_{it,3} + \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} x_{it,2} H_{i,s}^{1*} x_{is,3} \\ \vdots \\ 2 \sum_{t=1}^T H_{i,t}^{2*} x_{it,q-1} x_{it,q} + \sum_{t=1}^T \sum_{s \neq t} H_{i,t}^{1*} x_{it,q-1} H_{i,s}^{1*} x_{is,q} \end{pmatrix}. \quad (75) \end{aligned}$$

When $T = 1$, the score functions are as follow:

$$\mathbf{s}_\eta = \begin{pmatrix} s_\mu \\ s_\sigma \\ s_\beta \\ s_\gamma \end{pmatrix} = \begin{pmatrix} H_i^{1*} \\ H_i^{2*} \\ H_i^{1*} \mathbf{x}_i \\ H_i^{1*} \mathbf{z}_i \end{pmatrix}, \quad \begin{pmatrix} s_{\lambda_{\mu\mu}} \\ s_{\lambda_{\mu\sigma}} \\ s_{\lambda_{\sigma\sigma}} \\ s_{\lambda_{\mu\beta}} \\ s_{\lambda_{\sigma\beta}} \end{pmatrix} = \begin{pmatrix} H_i^{2*} \\ 3H_i^{3*} \\ 3H_i^{4*} \\ 2H_i^{2*} \mathbf{x}_i \\ 3H_i^{3*} \mathbf{x}_i \end{pmatrix}, \quad \text{and } s_{\lambda_{\beta\beta}} = \begin{pmatrix} H_i^{2*} x_{i,1}^2 \\ \vdots \\ H_i^{2*} x_{i,q}^2 \\ 2H_i^{2*} x_{i,1} x_{i,2} \\ \vdots \\ 2H_i^{2*} x_{i,1} x_{i,q} \\ 2H_i^{2*} x_{i,2} x_{i,3} \\ \vdots \\ 2H_i^{2*} x_{i,q-1} x_{i,q} \end{pmatrix}. \quad (76)$$

Notice that s_σ and $s_{\lambda_{\mu\mu}}$ are perfect collinear and, therefore, the Fisher information matrix associated with the proposed score function is singular under this reparameterization for data with $T = 1$.

B.3 Score function for testing $H_0 : m = M_0$ against $H_A : m = M_0 + 1$

The derivative of the reparameterized density w.r.t λ at ψ_τ^{h*} is zero similar to testing homogeneity case. With the constraint $\pi^{M_0} = 1 - \sum_{j=1}^{M_0-1} \pi^j$. The score functions s_{η_i} 's contain the first-order derivatives w.r.t π 's γ and ν at ψ_τ^{h*} :

$$\begin{aligned} \nabla_{\pi^j} l^h(\mathbf{w}; \psi_\tau^{h*}, \tau) &= \frac{f(\mathbf{w}; \gamma^*, \theta_0^{j*}) - f(\mathbf{w}; \gamma^*, \theta_0^{M_0^*})}{\sum_{j=1}^{M_0} \alpha_0^{j*} f(\mathbf{w}; \gamma^*, \theta_0^{j*})}; \\ \nabla_\gamma l^h(\mathbf{w}; \psi_\tau^{h*}, \tau) &= \frac{\sum_{j=1}^{M_0} \alpha_0^{j*} \nabla_\gamma f(\mathbf{w}; \gamma^*, \theta_0^{j*})}{\sum_{j=1}^{M_0} \alpha_0^{j*} f(\mathbf{w}; \gamma^*, \theta_0^{j*})}; \\ \nabla_\nu l^h(\mathbf{w}; \psi_\tau^{h*}, \tau) &= \frac{\nabla_\theta f(\mathbf{w}; \gamma^*, \theta_0^{h*})}{\sum_{j=1}^{M_0} \alpha_0^{j*} f(\mathbf{w}; \gamma^*, \theta_0^{j*})}. \end{aligned} \quad (77)$$

Define $H_{j,i,t}^{b*}$ as an abridged expression for $\frac{1}{b!} \frac{1}{\sigma_0^*} H^b \left(\frac{y_{it} - \mu_0^{j*} - x'_{it} \beta_0^{j*} - z'_{it} \gamma^*}{\sigma_0^{j*}} \right)$. Define the weight w_i^{j*} as

$$w_i^{j*} = \frac{\alpha_0^{j*} f(\{\mathbf{W}_{it}\}_{t=1}^T; \gamma^*, \theta_0^{j*})}{f_{M_0}(\{\mathbf{W}_{it}\}_{t=1}^T; \vartheta_{M_0}^*)}, \quad j = 1, \dots, M_0,$$

where $f_{M_0}(\{\mathbf{W}_{it}\}_{t=1}^T; \vartheta_{M_0}^*)$ is defined by equation (17).

As shown in section B.3, the score functions are:

$$\mathbf{s}_\alpha(\mathbf{w}_i) = \begin{pmatrix} \frac{f(\mathbf{w}|\theta_0^{1*}) - f(\mathbf{w}|\theta_0^{M_0*})}{\sum_t \alpha_0^{t*} f(\mathbf{w}|\theta_0^{t*})} \\ \vdots \\ \frac{f(\mathbf{w}|\theta_0^{M_0-1*}) - f(\mathbf{w}|\theta_0^{M_0*})}{\sum_t \alpha_0^{t*} f(\{\mathbf{W}_{it}^*\}_{t=1}^T | \theta_0^{t*})} \end{pmatrix}, \mathbf{s}_\mu(\mathbf{w}_i) = \begin{pmatrix} w_i^{1*} \sum_{t=1}^T H_{1,i,t}^{1*} \\ \vdots \\ w_i^{M_0*} \sum_{t=1}^T H_{M_0,i,t}^{1*} \end{pmatrix}, \mathbf{s}_\beta(\mathbf{w}_i) = \begin{pmatrix} w_i^{1*} \sum_{t=1}^T H_{1,i,t}^{1*} x_{it} \\ \vdots \\ w_i^{M_0*} \sum_{t=1}^T H_{M_0,i,t}^{1*} x_{it} \end{pmatrix},$$

$$\mathbf{s}_\sigma(\mathbf{w}_i) = \begin{pmatrix} w_i^{1*} \sum_{t=1}^T H_{1,i,t}^{2*} \\ \vdots \\ w_i^{M_0*} \sum_{t=1}^T H_{M_0,i,t}^{2*} \end{pmatrix}, \mathbf{s}_\gamma(\mathbf{w}_i) = \begin{pmatrix} w_i^{1*} \sum_{t=1}^T H_{1,i,t}^{1*} z_{it} \\ \vdots \\ w_i^{M_0*} \sum_{t=1}^T H_{M_0,i,t}^{1*} z_{it} \end{pmatrix}.$$

The score function for $\mathbf{s}_{\lambda\lambda}^h$ is obtained analogously to $\mathbf{s}_{\lambda\lambda}$ by replacing $H_{h,i,t}^{b*}$ with $H_{h,i,t}^{b*}$ for $b = 1, \dots, 4$ so that

$$\mathbf{s}_{\lambda_{\mu\sigma}}^h(\mathbf{w}_i) = w_i^{h*} \begin{pmatrix} \sum_{t=1}^T H_{h,i,t}^{2*} + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} H_{h,i,s}^{1*} \\ 3 \sum_{t=1}^T H_{h,i,t}^{4*} + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{2*} H_{h,i,s}^{2*} \\ 3 \sum_{t=1}^T H_{h,i,t}^{3*} + \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} H_{h,i,s}^{2*} \\ 2 \sum_{t=1}^T H_{h,i,t}^{2*} x_{it} + \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it} H_{h,i,s}^{1*} \\ 3 \sum_{t=1}^T H_{h,i,t}^{3*} x_{it} + 2 \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it} H_{h,i,s}^{2*} \end{pmatrix},$$

$$\mathbf{s}_{\lambda_\beta}^h(\mathbf{w}_i) = w_i^{h*} \begin{pmatrix} \sum_{t=1}^T H_{h,i,t}^{2*} x_{it,1}^2 + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it,1} H_{h,i,s}^{1*} x_{is,1} \\ \vdots \\ \sum_{t=1}^T H_{h,i,t}^{2*} x_{it,q}^2 + \frac{1}{2} \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it,q} H_{h,i,s}^{1*} x_{is,q} \\ 2 \sum_{t=1}^T H_{h,i,t}^{2*} x_{it,1} x_{it,2} + \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it,1} H_{h,i,s}^{1*} x_{is,2} \\ \vdots \\ 2 \sum_{t=1}^T H_{h,i,t}^{2*} x_{it,1} x_{it,q} + \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it,1} H_{h,i,s}^{1*} x_{is,q} \\ 2 \sum_{t=1}^T H_{h,i,t}^{2*} x_{it,2} x_{it,3} + \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it,2} H_{h,i,s}^{1*} x_{is,3} \\ \vdots \\ 2 \sum_{t=1}^T H_{h,i,t}^{2*} x_{it,q-1} x_{it,q} + \sum_{t=1}^T \sum_{s \neq t} H_{h,i,t}^{1*} x_{it,q-1} H_{h,i,s}^{1*} x_{is,q} \end{pmatrix}.$$

B.4 How to simulate the asymptotic distribution

C Other tables

Table 11: Parameter specification for null models with $M_0 = 1, 2, 3, 4$

	$M_0 = 1$
N	{100, 500}
T	{2, 5, 10}
a_n	(0.001, 0.005, 0.01, 0.05, 0.1, 0.2, 0.3, 0.4)
	$M_0 = 2$
N	{100, 500}
T	{2, 5, 10}
α	{(0.5, 0.5); (0.2, 0.8)}
μ	{(-1, 1), (-0.5, 0.5), (-0.8, 0.8)}
σ	{(1, 1), (1.5, 0.75), (0.8, 1.2)}
a_n	(0.01, 0.05, 0.1, 0.2, 0.3, 0.4)
	$M_0 = 3$
N	{100, 500}
T	{2, 10}
α	{(1/3, 1/3, 1/3); (0.25, 0.5, 0.25)}
μ	{(-4, 0, 4); (-4, 0, 5); (-5, 0, 5); (-4, 0, 6); (-5, 0, 6); (-6, 0, 6)}
σ	{(1, 1, 1); (0.75, 1.5, 0.75)}
a_n	(0.01, 0.05, 0.1, 0.2, 0.3, 0.4)
	$M_0 = 4$
N	{100, 500}
T	{2, 10}
α	{(0.25, 0.25, 0.25, 0.25)}
μ	{(-4, -1, 1, 4); (-5, -1, 1, 5); (-6, -2, 2, 6); (-6, -1, 2, 5); (-5, 0, 2, 4); (-6, 0, 2, 4)}
σ	{(1, 1, 1, 1); (1, 0.75, 0.5, 0.25)}
a_n	(0.01, 0.05, 0.1, 0.2, 0.3, 0.4)

Table 12: The estimated a_n -function based on the simulated nominal size

	<i>Dependent variable</i>			
	$\log\left(\frac{\hat{s}}{1-\hat{s}}\right) - \log\left(\frac{0.05}{1-0.05}\right)$			
	(1)	(2)	(3)	(4)
1/T	0.776*** (0.238)	-0.288*** (0.074)	0.611*** (0.050)	0.258*** (0.087)
1/N	28.143*** (10.127)	4.637 (3.124)	21.156*** (2.524)	8.585** (4.334)
$\log\left(\frac{a_n}{1-a_n}\right)$	-0.016 (0.019)	-0.101*** (0.009)	-0.111*** (0.007)	-0.128*** (0.030)
$\log\left(\frac{\omega(\vartheta_{M_0}; M_0)}{1-\omega(\vartheta_{M_0}; M_0)}\right)$		-0.197*** (0.029)	0.002 (0.006)	-0.013*** (0.003)
Constant	-0.616*** (0.113)	-0.811*** (0.047)	-0.680*** (0.060)	-0.735*** (0.068)
Observations	48	648	576	288

Note:

* p<0.1; ** p<0.05; *** p<0.01